

**CSES Working Paper Series**

**Paper # 37**

**Douglas Heckathorn**

**"Extensions of Respondent-Driven Sampling:**

**Analyzing Continuous Variables and Controlling for  
Differential Recruitment Using Dual-Component Sampling Weights**

**April 2007**

**EXTENSIONS OF RESPONDENT-DRIVEN SAMPLING:  
ANALYZING CONTINUOUS VARIABLES AND CONTROLLING FOR  
DIFFERENTIAL RECRUITMENT\***

(April 24, 2007)

Douglas D. Heckathorn  
Professor of Sociology  
Cornell University

- \* This research was made possible by grants from the *National Endowment for the Arts*, the *National Institutes of Health/National Institute on Drug Abuse*, and the *Centers for Disease Control and Prevention*. I thank Richard Campbell, Naihua Duan, Greg Duncan, Matthew Salganik, Michael Spiller III, Erik Volz, Cyprian Wejnert, and Carol Worthman for helpful comments and advice.

This is a preprint of an article accepted for publication in **Sociological Methodology** © 2007 American Sociological Association.

**EXTENSIONS OF RESPONDENT-DRIVEN SAMPLING:  
ANALYZING CONTINUOUS VARIABLES AND CONTROLLING FOR  
DIFFERENTIAL RECRUITMENT**

**ABSTRACT**

Respondent-driven sampling (RDS) is a network-based method for sampling hidden and hard-to-reach populations that has been shown to produce asymptotically unbiased population estimates when its assumptions are satisfied. This includes resolving a major concern regarding bias in chain-referral samples, that is, producing a population estimate that is independent of the seeds (initial subjects) with which sampling began. However, RDS estimates are limited to nominal variables, and one of the assumptions required for the proof of lack of bias is the absence of differential recruitment. One aim of this paper is to analyze the role of differential recruitment, quantify the bias it produces, and propose a new estimator that controls for it. The second aim is to extend RDS so that it can be employed to analyze continuous variables in a manner that controls for differential recruitment. Finally, the third aim is to describe means for carrying out multivariate analyses using RDS data. The analyses employ data from an RDS sample of 264 jazz musicians in the greater New York metropolitan area, taken in 2002.

## **EXTENSIONS OF RESPONDENT-DRIVEN SAMPLING:**

### **INTEGRATING RESPONDENT-DRIVEN AND MULTIPLICITY SAMPLING**

#### INTRODUCTION

Sampling what are termed “hard-to-reach” populations poses special problems because standard statistical sampling methods require a list of population members (i.e., a “sampling frame”) from which the sample can be drawn, and constructing the frame using methods such as household surveys is infeasible when the population is small relative to the general population and geographically dispersed, and when population membership involves stigma or the group has networks that are difficult for outsiders to penetrate (Sudman and Kalton 1986; Watters and Biernacki 1989; Spren 1992; Brown et al. 1999). Groups with these characteristics are relevant to research in many areas, including public health (e.g., drug users and commercial sex workers), public policy (e.g., illegal immigrants and the homeless), and arts and culture (e.g., jazz musicians and aging artists). In developing countries, inadequate public records compound sampling problems, and consequently much of the general population qualifies, in sampling terms, as “hidden.”

Sampling hidden populations has traditionally involved a dilemma. Some studies have employed probability sampling methods that provide incomplete coverage of the target population. For example, venue-based sampling (e.g., see MacKellar et al. 1996; Ramirez-Valles et al. 2005a) misses those who shun the large public venues, such as street corners and markets, from which subjects are recruited. Other studies have employed nonprobability sampling methods that provide more comprehensive coverage of the target population but yield only a convenience (i.e., nonstatistically valid) sample.

For example, snowball-type methods (Goodman 1961; Erickson 1979) start with a set of initial respondents (seeds), who refer their peers; these in turn refer their peers; and so on, as the sample expands from wave to wave. This approach has broader coverage because even those who shun public venues are reached through their social networks. Interest in these chain-referral methods has been fueled by recognition of this power to access members of hidden populations. As the literature on the “small world” asserts, even in a nation as large as the United States, every person is indirectly associated with every other person through approximately six intermediaries (Watts 2003). Therefore, everyone in the country could hypothetically be reached by the sixth wave of a maximally expansive chain-referral sample. However, inferences from a convenience sample cannot be made validly to the population from which the sample was drawn (Kalton 1983).

To overcome this dilemma, efforts have been made to transform snowball methods into probability sampling methods (Frank 1979; Snijders 1992; Spreen 1992; Frank and Snijders 1994). These are members of a relatively new class of probability sampling methods termed “adaptive” or “link-tracing” designs (Thompson and Frank 2000). This paper extends one such method, respondent-driven sampling (RDS) (Heckathorn 1997, 2002; Salganik and Heckathorn 2004, Volz and Heckathorn in press) in two ways. First, it introduces means for analyzing continuous variables that is based on delineating the relationship between RDS and a previously introduced method for analyzing chain-referral data, Sirken’s (1970) multiplicity sampling. Second, based on partitioning the RDS sampling weight into a multiplicity-based component and a component based on analyzing cross-group recruitment patterns, it introduces means for controlling bias due to differential recruitment, in which subgroups have both differing

recruitment patterns and recruitment effectiveness, and hence the more effective recruiting group's recruitment patterns differentially affect sample composition.

Section I reviews the fundamentals of multiplicity sampling and RDS and introduces the required notation. Section II shows how the weight for a dichotomous variable can be partitioned into multiplicity and recruitment components, and introduces means for weighting continuous variables based on the partitioning of the sampling weight. Section III specifies the conditions under which differential recruitment can introduce bias into the RDS population estimator, and introduces a new estimator that controls for that bias. Finally, the conclusion discusses potential areas of further refinement of the RDS method.

## **I. Respondent-Driven Sampling and its Relationship to Multiplicity Sampling: Basic Concepts**

### **Multiplicity Sampling**

Multiplicity sampling was developed by Sirken (1970) in the late 1960s for sampling rare populations. The approach is straightforward. A multiplicity survey differs from a conventional survey because each case may appear more than once. For example, a telephone directory may have multiple entries for the same household. Consequently, when households are sampled from such a directory, they must be weighted based on multiplicity; for example, those with three phones have a weight only one-third that of households with a single listed phone. This approach is useful for increasing the efficiency with which household surveys can estimate the prevalence of rare events. The respondent is asked not only whether a condition affects his or her household, but also whether it affects a specified group of other households, such as those of surviving

children and siblings. In this way, information regarding the event becomes available not only from the household surveyed, but also from other households to which it is connected. Multiplicity arises because an event can be reported from multiple sources. In this way, the ability to detect rare events is increased.

The multiplicity approach was extended to snowball samples by Rothbart, Fine, and Sudman (1982). They proposed adding to the survey a question regarding the number of eligible respondents known to the respondent. The size of this network then provides the basis for a multiplicity adjustment, in which respondents are weighted by the reciprocal of their network sizes. The intuition underlying this approach is that respondents with large networks will have a greater probability of inclusion, because more recruitment paths lead to them, and thus respondents must be weighted by the reciprocal of their network sizes; that is, for any individual  $i$ , its weight is  $1/D^i$ . The multiplicity weight for any individual  $i$ ,  $MW^i$ , can therefore be defined as

$$(1) \quad MW^i = \frac{1}{D^i}$$

(Here and elsewhere in the paper, *superscripts* are employed to index *individuals*, and *subscripts* are employed to index *groups*.) These weights can then be employed to analyze any variable; for example, Table I shows the multiplicity estimates for two nominal variables, gender and having received airplay; Table II shows the multiplicity estimates for two continuous variables divided by quintiles, age and degree; and Figure 1 shows them analyzed as continuous variables.

TABLE I ABOUT HERE

TABLE II ABOUT HERE

FIGURE 1 ABOUT HERE

A limitation of this approach is that it fails to control for bias resulting from differential recruitment. For example, among New York City jazz musicians, recruitment effectiveness varied by gender; though females made up 26% (65/243) of respondents, they produced 37% (91/243) of the recruits (see Table IA). Consequently, the female patterns of recruitment can be expected to have differentially affected the sample. This is consequential, because recruitment patterns differed by gender; females recruited 44% (40/91) other females whereas males recruited only 16% females (25/152), so females were oversampled. Therefore, both elements required for the presence of differential recruitment bias are present—differential recruitment effectiveness and different recruitment patterns. Based on the same logic, musicians who received airplay were undersampled.

### **Respondent-Driven Sampling**

RDS employs both the degree data upon which multiplicity sampling depends, and also information on patterns of recruitment within the sample, specifically, the proportions of recruitment across groups. The RDS estimator is derived from an analysis of network structure. (For comprehensive descriptions, see Heckathorn 2002, Salganik and Heckathorn 2004, and Volz and Heckathorn in press). The connection between network structure and a population estimator is based on the *reciprocity model*, named for a feature of the networks of friends and acquaintances through which peer recruitment takes place. Ties are reciprocal, so a link from any individual  $i$  to  $j$  implies that a link also exists from  $j$  to  $i$ . Consequently, no distinction need be made between ties to an individual (the in-degree) and ties from the individual to others (the out-degree), since the two are equivalent. In such systems, reciprocity also extends to ties linking groups.



The RDS estimator is based on this elemental feature of reciprocal networks. In a two-group system, the number of ties from group X to Y is the product of four parameters. The first is the population size, N. For example, in a system consisting of two disjoint groups, X and Y, with union equal to the population, and where  $N_X$  is the number of Xs, and  $N_Y$  the number of Ys, therefore,

$$(2) \quad N = N_X + N_Y$$

The second is the proportional size of the group,  $P_X$ , i.e.,

$$(3) \quad P_X = \frac{N_X}{N}$$

The third is the group's mean degree. Where  $T_X$  is the number ties of group X, the group's mean degree,  $D_X$ , is

$$(4) \quad D_X = \frac{T_X}{N_X}$$

The model therefore relies on the most basic measure of centrality. This measure is used rather than more complex measures of centrality, such as eigenvector centrality (Bonacich, 1972), that reflect the relative influence of individuals or groups. What is relevant for this model is merely the number of connections to other nodes, because this reflects the size of the pool from which potential recruits are drawn. The fourth and final parameter is the proportion of cross-cutting ties. Where  $T_{XY}$  is the number of ties from X to Y, the group's proportion of cross ties,  $S_{XY}$ , is

$$(5) \quad S_{XY} = \frac{T_{XY}}{T_X}$$

The number of cross-group ties from group X to Y,  $T_{XY}$ , is the product of these four terms. That is,

$$(6) \quad T_{XY} = N P_X D_X S_{XY}$$

In this expression, the product of the first two terms is the size of the group (i.e.,  $N * P_X = N_X$ ), the product of the first three terms is the number of ties of the group (i.e.,  $N_X * D_X = T_X$ ), and consequently the product of all four terms is the number of cross-cutting ties (i.e.,  $T_X * S_{XY} = T_{XY}$ ).

When the ties in a system are reciprocal, such that in-degrees and out-degrees are equivalent, the number of cross-cutting ties will be equal in each direction, i.e., for groups X and Y,

$$(7) \quad T_{XY} = T_{YX}$$

Given that the groups' proportional sizes sum to one, and expanding the expression for cross-cutting ties in each direction yields the following equation system,

$$(8) \quad \begin{aligned} 1 &= P_X + P_Y \\ N P_X D_X S_{XY} &= N P_Y D_Y S_{YX} \end{aligned}$$

These can be solved to yield group X's proportional size,  $P_X$ , as follows:

$$(9) \quad P_X = \frac{S_{YX} D_Y}{S_{YX} D_Y + S_{XY} D_X}$$

This equation provides the basis for an estimator for proportional group size,  $\widehat{P}_X$ , based on two types of network information. One is the estimated proportion of cross-cutting ties (the "S" terms), and estimated mean network size (the "D" terms), both of which, as will be seen below, can be derived from chain-referral data.<sup>1</sup> That is,

$$(10) \quad \widehat{P}_X = \frac{\widehat{S}_{YX} \widehat{D}_Y}{\widehat{S}_{YX} \widehat{D}_Y + \widehat{S}_{XY} \widehat{D}_X}$$

For example, using the estimates for degree and cross-cutting ties from Table IA, the estimated proportion of females,  $\widehat{P}_F$ , is calculated as follows:

$$(11) \quad \widehat{P}_F = \frac{0.164 \cdot 109.255}{0.164 \cdot 109.255 + 0.56 \cdot 102.566} = 0.238$$

This estimator contrasts with the “face value” estimator typically used in analyzing chain-referral data, the proportion of each group in the sample, i.e., where  $n_X$  is the number of Xs in the sample, and  $n$  is the sample size, the sample proportion,  $C_X$ , is

$$(12) \quad C_X = \frac{n_X}{n}$$

The RDS estimator has been shown to be asymptotically unbiased (Salganik and Heckathorn 2004), which means that bias is on the order of  $1/[\text{sample size}]$ , so bias is negligible in samples of meaningful size (Cochran 1977). The proof is based on six assumptions about the sampling process:

- (1) Respondents know one another as members of the target population, so ties are reciprocal.
- (2) Respondents are linked by a network composed of a single component.
- (3) Sampling occurs with replacement.
- (4) Respondents can accurately report their personal network size, defined as the number of relatives, friends, and acquaintances who fall within the target population.
- (5) Peer recruitment is a random selection from the recruiter’s network.
- (6) Each respondent recruits a single peer.

The first three assumptions serve to specify the conditions under which RDS is an appropriate sampling method. First, peer recruitment is a feasible sampling strategy only

if respondents know one another as members of the target population. Consequently, it would not be suitable for sampling tax cheats, who can be friends and not know they share membership in that hidden population. However, it is suitable for sampling populations linked by a “contact pattern,” such as reciprocal ties created when jazz musicians perform with one another or when drug users purchase drugs. Second, ties must be dense enough to sustain the chain-referral process. When respondents recruit friends and acquaintances, this is rarely problematic, because populations linked by a contact pattern tend to be gregarious. For example, Heckathorn and Jeffri (2003) found that the typical New York City jazz musician knew about 100 other musicians and none knew fewer than 20, a number greater than that generally required for a network to form a single large component. In contrast, allowing recruitment only of musicians who perform together would not be advisable, because the network would comprise many small components. Third, sampling is assumed to occur with replacement, so recruitments do not deplete the set of respondents available for future recruitment. The implication is that the sampling fraction should be small enough for a sampling-with-replacement model to be appropriate.

The fourth assumption states that respondents can accurately report the number of relatives, friends, and acquaintances who are members of the target population. Studies of the reliability of network indicators suggest that this is one of the more reliable indicators (Marsden 1990). Furthermore, the RDS population estimator depends not on absolute but on relative degree, so variations in name generators that inflate or deflate the reports in a linear manner have no effect on the estimates. However, violations of this assumption about accurate reporting remain a source of bias on which additional research would be

useful.

The fifth assumption specifies that respondents recruit as though they are choosing randomly from their networks. This is based on the assumption that recruitment will be nonbiased because respondents would lack an incentive or ability to coordinate to selectively recruit any particular group. Evidence for this assumption has been provided by studies that compared self-reported network composition with actual recruitment behavior and found a strong association (Heckathorn et al. 2002; Wang et al. 2005), and also by a study in which an “index of reciprocity” measured the fit between the reciprocity model and recruitment patterns (Ramirez-Valles et al. 2005b).

The plausibility of the random recruitment assumption is determined, in part, by the research design. For example, if a research site were located in a dilapidated building in a high-crime neighborhood, recruiting residents of the neighborhood might be easy, but recruiting peers from more comfortable neighborhoods who felt threatened in such neighborhoods might prove difficult, so sampling would be non-random. However, if research identifies neutral turf in which all potential respondents feel safe, the random recruitment assumption is made more plausible. Similarly, if incentives are offered that are salient to respondents from all income groups (e.g., a choice between receiving a monetary reward and making a contribution to a charity of the respondent’s choice), the random recruitment assumption is made more plausible. Also, if respondents who live in areas distant from the interview site have limited access to means of transportation, that group will be under sampled. If either additional interview sites are made available that are closer to them, or if the boundaries of the target population are reduced to include only those respondents with ready access to the original interview site, the random

recruitment assumption is made more plausible. Finally, if one formulated the network-size question using a time frame of five years, it is probable that persons who had been seen between one and five years ago would be under sampled, so sampling would be non-random. However, if research shows that the great majority of recruits are seen at least once per month, using the time frame of a month would make the random recruitment assumption more plausible.

The sixth assumption, that each respondent recruits exactly one peer, serves to preclude differential recruitment. This is an especially problematic assumption because some respondents fail to recruit, so the chain-referral process can die out. For example, in Klovdahl's (1989) "random walk" approach, recruitment is limited to three waves, and only one-quarter of chains attain that length. The sample therefore consists of multiple short, linear chains. This introduces the potential for differential recruitment, e.g., groups that recruit less effectively than others will be overrepresented on the recruitment chain's terminal node. In RDS studies it is customary to establish a nonunitary quota of permitted recruitments, generally a limit of three or four. This number has been found to produce robust referral chains. For example, in a study of New York City drug users the quota was three, recruitment began with 8 seeds and over the course of 18 waves yielded a sample of 618 (Abdul-Quader 2006). The NYC jazz study employed a quota of four, with 10 seeds, and one seed produced a recruitment chain with more than 100 other respondents over the course of 10 waves. Such recruitment introduces considerable potential for differential recruitment, and to varying degrees this occurs in all RDS data sets. An aim of this paper is to propose a new estimator that controls for this source of bias.

The RDS estimator is calculated from two distinct terms, the proportion of cross-cutting ties between groups, and the mean degree of each group. The following sections discuss the ways in which these parameters can be estimated based on chain-referral data.

*Estimating the Proportion of Cross-Cutting Ties*

Deriving the first type of information requires documenting who recruited whom, usually based on recruitment coupons with unique serial numbers that are recorded when given to the recruiter and again when returned by the recruit. For each variable to be analyzed, a recruitment matrix,  $R$ , is calculated, where  $R_{XY}$  is the number of recruitments by members of group X of members of group Y,

$$(13) \quad R = \begin{bmatrix} R_{XX} & R_{XY} \\ R_{YX} & R_{YY} \end{bmatrix}$$

In this matrix, the row sums reflect the number of recruitments *by* members of each group, e.g., the number of recruitments by group X,  $RB_X$ , is

$$(14) \quad RB_X = R_{XX} + R_{XY}$$

For example, in the analysis of recruitment by gender (see Table IA), 152 respondents were recruited by males, consisting of 127 males and 25 females, and 91 respondents were recruited by females, consisting of 51 males and 40 females.

In the recruitment matrix, the column sums reflect the number of recruitment *of* members of each group, e.g., the number of recruitments of group X,  $RO_X$ , is

$$(15) \quad RO_X = R_{XX} + R_{YX}$$

For example, in the analysis of recruitment by gender (see Table IA), there were 178 recruitments of males, 127 by other males and 51 by females.

The number of cases in the recruitment matrix, it should be noted, is necessarily

less than the sample size, because seeds do not have a recruiter. In the absence of missing data, the number of cases in the recruitment matrix, (i.e., the number of recruits, RO, or equivalently, the number of recruitments, RB), is the sample size, n, less the number of seeds,  $n_s$ ,

$$(16) \quad RO = RB = n - n_s$$

Hence for purposes of RDS analyses, the effective sample size is the number of cases less the number of seeds. Of course, this number is further reduced by missing data. In RDS missing data is especially problematic, because when data for a respondent is missing, neither its recruitment, nor recruitments by it, appear in the recruitment matrix. For example, if A is recruited by B, who recruits C, D, and E, and B's data is missing, then the number of recruitments is reduced by four, because recruitments  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $B \rightarrow D$ , and  $B \rightarrow E$ , are lost.

Based on the recruitment matrix, the recruitment selection proportions can be calculated. These are terms that serve as the estimators for the proportion of cross-group ties. Specifically, the ratio of cross-group recruitments,  $R_{XY}$ , and total recruitments by the group,  $RB_X$  provides an estimator for  $S_{XY}$ , that is,

$$(17) \quad \widehat{S}_{XY} = \frac{R_{XY}}{RB_X}$$

For example, in Table IA's analysis of recruitment by gender, the estimated selection proportion from females to males is  $51/91 = 0.56$ .

This estimator has been shown to be unbiased (Salganik and Heckathorn 2004:214), because based on the random-recruitment assumption the proportion of ties that become the basis for peer recruitment must be equal across subgroups. That is, if the sampling fraction for X's ties is SF, then the number of recruitments by X is the product



of X's number of ties,  $T_X$ , and the sampling fraction, SF, i.e.,

$$(18) \quad RB_X = T_X \cdot SF$$

Furthermore, from the random recruitment assumption, the sampling fraction for X's cross cutting ties,  $T_{XY}$ , must be the same, for otherwise cross-cutting ties would be either over or undersampled, hence

$$(19) \quad R_{XY} = T_{XY} \cdot SF$$

Therefore, the above expression for  $\widehat{S}_{XY}$  can be expanded as follows,

$$(20) \quad \widehat{S}_{XY} = \frac{R_{XY}}{RB_X} = \frac{T_{XY} SF}{T_X SF}$$

Given that the SF terms cancel,  $R_{XY}/RB_X$  provides an unbiased estimator for  $T_{XY}/T_X$ .

The implication is that the first element from which the RDS estimator is calculated, the cross-group recruitment proportion, is free from bias due to differential recruitment.

### *Estimating Degree*

The second element from which the RDS estimator is calculated is the estimated mean degree of each group. This estimator employs a multiplicity approach. That is, consistent with this approach, respondents are assumed to be recruited in proportion to their degree. Respondents of higher degree are oversampled, so in estimating a group's degree, respondents are weighted by the inverse of their degree (Salganik and Heckathorn 2004:215, Volz and Heckathorn, in press). For any group X, where  $n_X$  is the number of respondents falling within that group, and  $D^i$  is the degree of respondent i, the estimated mean network size for group X,  $\widehat{D}_X$ , is

$$(21) \quad \widehat{D}_x = \frac{n_X}{\sum_{i=1}^{n_X} \frac{1}{D^i}}$$

Salganik and Heckathorn (2004:218) showed that both the numerator and the denominator of this expression correspond to Hansen-Hurwitz (1943) estimators, which are known to be unbiased (Brewer and Haif 1983). It is also known (Cochran 1977) that the ratio of these estimators is asymptotically unbiased, with bias on the order of  $1/[\text{sample size}]$ , so bias falls as sample size increases.

The RDS estimator (equation 10) includes degree estimates in both the numerator and the denominator, each of which is asymptotically unbiased. The ratio of asymptotically unbiased estimators is also asymptotically unbiased; therefore the RDS estimator is also asymptotically unbiased (Heckathorn and Salganik 2004:219).

A limitation of this approach is that it cannot be used to analyze continuous variables. For example, respondents ranged in age from 20 to 101, and when partitioned by year, there were 54 distinct ages. Yet analysis based on a  $54 \times 54$  matrix with 2,916 cells among which the 264 respondents would be distributed is infeasible. Of course, the sample could be aggregated, e.g., divided by quartiles or quintiles, but this would entail loss of information. Section II introduces means for analyzing continuous variables that reduces but does not wholly eliminate this loss of information.

A second limitation of the approach is that this way of estimating mean degree does not control for differential recruitment. For example, if respondents of high degree associate differentially with one another and also recruit more effectively than those of lower degree, high-degree respondents will be oversampled. Section III introduces means for controlling for this source of bias based on Section II's reformulation of the RDS sampling weight to accommodate analysis of continuous variables.

## **II. Dual-Component Sampling Weights**

### *Partitioning the RDS Sampling Weight*

A step toward endowing RDS with the multiplicity approach's ability to analyze continuous variables is to divide the RDS sampling weight into two components, one that adjusts for differential recruitment, and one that adjusts for differences in degree. The RDS sampling weight for any group  $X$ ,  $W_x$ , is the ratio of the population estimate for the group,  $\widehat{P}_x$ , and the proportional composition of the sample,  $C_x$ . Therefore,

$$(22) \quad W_x = \frac{\widehat{P}_x}{C_x}$$

Dividing this weight into recruitment and degree components is based on separating the effects of differential recruitment and differences in degree. This can be done by projecting what the sample composition would have been in the absence of both factors. Heckathorn (2002) suggests a means by which this can be done. It involves modeling the recruitment process as a first-order Markov process. The state space is fixed, with each group corresponding to a state, and the recruitment proportions in the recruitment matrix are interpreted as transition probabilities. The sampling process is then modeled as sequences of states governed by the transition probabilities. For example, if sampling began (wave zero) with a female seed, from Table IA, there would be a 44% probability that the next respondent would be female, and a 56% probability that the next recruit would be male. If the first-wave recruit was male, there would be a 16% probability that the next respondent would be female, and an 84% probability that the next recruit would be male. The sample expands in this stochastic manner in subsequent waves. When modeled in this manner, the sample reaches an equilibrium composition that is independent of the state (i.e., initial respondent, or equivalently, the "seed") from which it

began (Kemeny and Snell 1960; Heckathorn 1997). In a two-state system, with groups X and Y, the equilibrium is defined by the following equation system,

$$(23) \quad \begin{aligned} 1 &= E_X + E_Y \\ E_X &= S_{XX}E_X + S_{YX}E_Y \end{aligned}$$

Substituting  $1 - S_{XY}$  for  $S_{XX}$ , and solving for  $E_X$  yields the following,

$$(24) \quad E_X = \frac{S_{YX}}{S_{YX} + S_{XY}}$$

From this expression, it is clear that equilibrium is a term that is meaningful only at the group level, for it is defined not based on individual or unitary group attributes, but rather based on the proportions of cross-cutting ties across groups.

The equilibrium provides the means to project what the sample composition would have been in the absence of differences in degree (Heckathorn 2002:25). This can be shown by assuming that both groups have equal degree and calculating the RDS estimator; that is, if groups X and Y both have equal degree, then  $\widehat{D}_X$  can be substituted for  $\widehat{D}_Y$  in equation 10's expression for the RDS estimator, and by algebraic manipulation, this expression can be simplified as follows:

$$(25) \quad \widehat{P}_X = \frac{\widehat{S}_{YX}}{\widehat{S}_{YX} + \widehat{S}_{XY}} \quad \text{if } \widehat{D}_X = \widehat{D}_Y$$

This equation is equivalent to that for the Markov equilibrium. Consequently, the equilibrium can be seen as projecting what the sample composition would have been had degrees been uniform across groups.

A similar argument (Heckathorn 2002:21) shows that the Markov equilibrium also projects what the sample composition would have been in the absence of differential recruitment. The equilibrium is calculated exclusively from the transition probabilities,

and above it was shown that these are independent of differential recruitment.

Consequently, a term calculated from the transition probabilities is also independent of differential recruitment.

Because the Markov equilibrium provides a baseline indicator showing what the sample composition would have been in the absence of both differential recruitment and differences in degree, it thereby provides the means for partitioning the RDS sampling weight to disentangle these two factors, as follows:

$$(26) \quad W_X = \frac{\widehat{P}_X}{C_X} = \frac{\widehat{P}_X}{\widehat{E}_X} \cdot \frac{\widehat{E}_X}{C_X}$$

Here the term  $\widehat{P}_X / \widehat{E}_X$  can be termed the degree component,  $DC_X$ ,

$$(27) \quad DC_X = \frac{\widehat{P}_X}{\widehat{E}_X}$$

When degrees are equal,  $\widehat{P}_X = \widehat{E}_X$ , so their ratio has the neutral value of unity, i.e.,  $DC_X$

= 1. In contrast, if X has greater mean degree than Y, group X is oversampled, so the estimated proportional size of the group must be correspondingly deflated such that

$\widehat{P}_X < \widehat{E}_X$ , and the value of  $DC_X$  is then less than one. By the same logic, if group X has a

smaller mean degree than Y, then  $\widehat{P}_X > \widehat{E}_X$ , and the value of  $DC_X$  is greater than one.

This inverse relationship between DC and the group's mean degree derives from the presence of the latter in the denominator of the population estimator.

The second term of the partitioned weight,  $\widehat{E}_X / C_X$ , can be termed the recruitment component,  $RC_X$ , because in the absence of differential recruitment, the two are equal, so this term has the neutral value of unity; hence

$$(28) \quad RC_X = \frac{\widehat{E}_X}{C_X}$$

As thus defined, the product of the degree and recruitment components yields the sampling weight,

$$(29) \quad W_X = DC_X \cdot RC_X$$

In the definition of the degree and recruitment components, it is useful to be clear about the role played by this equilibrium. The question is not whether it is behaviorally plausible to model recruitment as a first-order Markov process, though evidence for this has been presented (Heckathorn 1997:83), but rather that this term provides the means for abstracting from the effects of both differential recruitment and differences in degree. For example, the degree term does not appear in the equation for the equilibrium; it is calculated exclusively from the selection proportions. Consequently, for the equilibrium degrees are irrelevant; and comparing the value of that term with the population estimate for which degrees matter provides the means for quantifying the effects that degrees have on the estimation process.

When viewed in light of the distinction between the degree and recruitment components, the gender and airplay variables represent contrasting cases. For gender, the recruitment component is the principal determinant of the sample weight. The mean departure from unity of the recruitment component for males and females is 2.9 times greater than that of their degree component. For airplay the relationship is reversed. The mean departure from unity of the degree component for musicians with and without airplay is 12.9 times greater than that of their recruitment component. This greater dependence of airplay on the degree component reflects the dependence of degree on

airplay; musicians with airplay had networks that were 47% larger than those without airplay. In contrast, gender is a weak determinant of degree: males have only 5.5% larger networks than females. (For illustrations of the above-described calculation procedures, see Appendix A.)

*Extending the Analysis to Nondichotomous Variables*

A second step in extending the dual-component analysis to continuous variables is showing how three-category and larger variables are analyzed. The Markov equilibrium extends to non-dichotomous variables in a straightforward manner. For a system with M categories, calculating the equilibrium requires solving the following system of equations (Kemeny and Snell 1960),

$$\begin{aligned}
 & 1 = \widehat{E}_1 + \widehat{E}_2 + \dots + \widehat{E}_M \\
 & \widehat{E}_1 = \widehat{S}_{11}\widehat{E}_1 + \widehat{S}_{21}\widehat{E}_2 + \dots + \widehat{S}_{M1}\widehat{E}_M \\
 (30) \quad & \widehat{E}_2 = \widehat{S}_{12}\widehat{E}_1 + \widehat{S}_{22}\widehat{E}_2 + \dots + \widehat{S}_{M2}\widehat{E}_M \\
 & \quad \quad \quad \vdots \\
 & \widehat{E}_{(M-1)} = \widehat{S}_{1(M-1)}\widehat{E}_1 + \widehat{S}_{2(M-1)}\widehat{E}_2 + \dots + \widehat{S}_{M(M-1)}\widehat{E}_M
 \end{aligned}$$

This consists of a system of M linear equations, with M unknowns so when the selection proportions are known, the equilibrium has a unique solution.

Calculating the RDS population estimator for variables with more than three categories is less straightforward, though it involves solving a somewhat similar system of equations. As in the two categories case (equation 8), the first equation states that proportional population sizes must sum to one. The other equations express the reciprocity principle for each of the  $M*(M-1)/2$  pairs of groups. For example, a system with three disjoint groups is described by four equations, as follows, where the population size parameter, N, is omitted because it cancels out,

$$\begin{aligned}
& 1 = \widehat{P}_1 + \widehat{P}_2 + \widehat{P}_3 \\
(31) \quad & \widehat{P}_1 \widehat{D}_1 \widehat{S}_{12} = \widehat{P}_2 \widehat{D}_2 \widehat{S}_{21} \\
& \widehat{P}_1 \widehat{D}_1 \widehat{S}_{13} = \widehat{P}_3 \widehat{D}_3 \widehat{S}_{31} \\
& \widehat{P}_2 \widehat{D}_2 \widehat{S}_{23} = \widehat{P}_3 \widehat{D}_3 \widehat{S}_{32}
\end{aligned}$$

When the “D” and “S” terms are calculated in the above-described manner, this yields a system of four linear equations with three unknowns, the “P” terms. Consequently, the system is over-determined, because the number of equations exceeds the number of unknowns. This issue arises for any variable that contains three or more categories. The most standard statistical approach to solving such systems is linear least squares (Farebrother 1988), which employs a regression-like logic to reconcile conflicts among the equations. For a discussion of this approach, as applied to RDS, see Heckathorn (2002:23).

An alternative approach, termed data smoothing (Heckathorn 2002:24-25), derives from drawing information regarding the population from the reciprocity model. The essential idea is that if ties in the system are reciprocal, if all groups recruit with equal effectiveness (i.e., for any group X,  $RO_X = RB_X$ ), and recruitments from personal networks are random, then cross-group recruitments will be equal for each pair of groups (i.e., for any groups X and Y,  $R_{XY} = R_{YX}$ ). Consequently, any differences reflect merely stochastic variation in the recruitment process, so the best estimate for the number of cross-recruitments between each pair of groups is the mean of recruitments in each direction. This form of data reduction has several advantages. First, by reducing the number of terms from which population estimates are calculated, it solves the problem of over-determination because the additional equations that produce the over-determination problem are rendered redundant and hence can be ignored. Second, each cross-



recruitment term is calculated from twice as much data, so estimates based upon them become more efficient, as reflected in narrower confidence intervals (Volz and Heckathorn, in press). For example, if Table IIA’s analysis of age is carried out using linear-least squares, the design effect is 2.1, but this falls to 1.6 when data smoothing is used. Finally, data smoothing preserves a feature crucial for the dual-component approach, in which the RDS population estimator equals the equilibrium when degrees are equal. For this reason, data smoothing will be employed in this paper for all three-category and larger variables. It will not be employed for dichotomous variables, for point estimates would be unaffected, however for purposes of variance estimation data smoothing is useful even in the two-category case (see Volz and Heckathorn, in press).

Data smoothing is a two-step process. One step is projecting what the recruitment matrix would have looked like in the absence of differential recruitment. This requires transforming the matrix using two conditions: (1) recruitment patterns, as reflected in the selection proportions, do not change; (2) the row and column sums are equal, so recruitment effectiveness is equal for all groups (i.e., for any group X,  $RO_X=RB_X$ ). This transformation has been termed “demographic adjustment” (Heckathorn 2002:21, Volz and Heckathorn, in press) and when used in other contexts is termed “raking.” Each element in the transformed recruitment matrix is the product of three terms, the selection proportion, the equilibrium for the recruiter’s group, and the total number recruitments. For example, for recruitment count,  $R_{XY}$ , the corresponding demographically adjusted recruitment count,  $R^*_{XY}$  is  $\widehat{S}_{XY} \widehat{E}_X RB$ . For a system with M categories, the adjusted recruitment matrix,  $R^*$  is,

$$(32) \quad R^* = \begin{bmatrix} \widehat{S}_{11}\widehat{E}_1RB & \widehat{S}_{12}\widehat{E}_1RB & \cdots & \widehat{S}_{1M}\widehat{E}_1RB \\ \widehat{S}_{21}\widehat{E}_2RB & \widehat{S}_{22}\widehat{E}_2RB & \cdots & \widehat{S}_{2M}\widehat{E}_2RB \\ \vdots & \vdots & \ddots & \vdots \\ \widehat{S}_{M1}\widehat{E}_M RB & \widehat{S}_{M2}\widehat{E}_M RB & \cdots & \widehat{S}_{MM}\widehat{E}_M RB \end{bmatrix} = \begin{bmatrix} R_{11}^* & R_{12}^* & \cdots & R_{1M}^* \\ R_{21}^* & R_{22}^* & \cdots & R_{2M}^* \\ \vdots & \vdots & \ddots & \vdots \\ R_{M1}^* & R_{M2}^* & \cdots & R_{MM}^* \end{bmatrix}$$

By inspection, it is clear that this transformation of the recruitment matrix does not alter the selection proportions, because each is multiplied by the same constant, e.g.,  $\widehat{E}_1RB$  in the first row. Consequently, the transformation satisfies the first condition. The second condition is satisfied because the inclusion of the equilibrium in each term ensures that the row and column proportions will each equal that equilibrium, and hence the corresponding row and column counts will be equal. When transformed in this manner, the matrix has a simple structure. Not only are the row and column sums equal, but they also equal the equilibrium, i.e., for any group  $X$ ,  $RB_X^* = RO_X^* = \widehat{E}_X$ .

If the assumptions of the RDS model are satisfied, such that ties are reciprocal and recruitment is random, then differences in cross-recruitment counts reflect only stochastic variation. Consequently, the matrix can be smoothed by taking the mean of these counts, to yield a smoothed recruitment matrix,  $R^{**}$  as follows,

$$(33) \quad R^{**} = \begin{bmatrix} \widehat{S}_{11}\widehat{E}_1RB & \frac{(\widehat{S}_{12}\widehat{E}_1RB) + (\widehat{S}_{21}\widehat{E}_2RB)}{2} & \cdots & \frac{(\widehat{S}_{1M}\widehat{E}_1RB) + (\widehat{S}_{M1}\widehat{E}_M RB)}{2} \\ \frac{(\widehat{S}_{12}\widehat{E}_1RB) + (\widehat{S}_{21}\widehat{E}_2RB)}{2} & \widehat{S}_{22}\widehat{E}_2RB & \cdots & \frac{(\widehat{S}_{2M}\widehat{E}_2RB) + (\widehat{S}_{M2}\widehat{E}_M RB)}{2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{(\widehat{S}_{1M}\widehat{E}_1RB) + (\widehat{S}_{M1}\widehat{E}_M RB)}{2} & \frac{(\widehat{S}_{2M}\widehat{E}_2RB) + (\widehat{S}_{M2}\widehat{E}_M RB)}{2} & \cdots & \widehat{S}_{MM}\widehat{E}_M RB \end{bmatrix} = \begin{bmatrix} R_{11}^{**} & R_{12}^{**} & \cdots & R_{1M}^{**} \\ R_{21}^{**} & R_{22}^{**} & \cdots & R_{2M}^{**} \\ \vdots & \vdots & \ddots & \vdots \\ R_{M1}^{**} & R_{M2}^{**} & \cdots & R_{MM}^{**} \end{bmatrix}$$

The effect of this transformation is to render the recruitment matrix “reciprocity compatible” (Heckathorn 2002:32) in the sense that the additional equations that produce the problem of over determination in three-category and larger systems is resolved,

because the excess equations become redundant.

After data smoothing is complete, the smoothed recruitment matrix becomes the basis for all calculations, so all terms dependent on that matrix must be recalculated, including the row sums, the column sums, the selection proportions, and the equilibrium. For example, based on the smoothed selection proportions, and the estimated degrees (which are not altered by data smoothing), the smoothed population estimate is calculated as follows in a system with M groups,

$$\begin{aligned}
 (34) \quad & 1 = \widehat{P}_1^{**} + \widehat{P}_2^{**} + \widehat{P}_3^{**} + \cdots + \widehat{P}_M^{**} \\
 & \widehat{P}_1^{**} \widehat{D}_1 \widehat{S}_{12}^{**} = \widehat{P}_2^{**} \widehat{D}_2 \widehat{S}_{21}^{**} \\
 & \widehat{P}_1^{**} \widehat{D}_1 \widehat{S}_{13}^{**} = \widehat{P}_3^{**} \widehat{D}_3 \widehat{S}_{31}^{**} \\
 & \quad \quad \quad \vdots \\
 & \widehat{P}_1^{**} \widehat{D}_1 \widehat{S}_{1M}^{**} = \widehat{P}_M^{**} \widehat{D}_M \widehat{S}_{M1}^{**}
 \end{aligned}$$

This expression provides the means for calculating population estimates for a system with M categories by solving a system of M equations with M unknowns. A parallel with the n-category expression for the Markov equilibrium is present, because if degrees are equal across groups, the degree term drops out, thereby generalizing to the n-category case the conclusion that equal degrees imply equality between the population estimate and the equilibrium. In the balance of this paper, smoothing will be employed when calculating estimators for all three-category and larger variables, however, to simplify the notation the double asterisk indicating smoothing will not be shown. (For an example showing how to calculate a smoothed population estimate, see Appendix B.)

#### *Calculating Individualized Degree Components of the RDS Sampling Weight*

The next step in extending RDS to permit analysis of continuous variables is to devise means for calculating the degree component in an individualized manner. It is

inherent in chain-referral samples, including RDS, that respondents of high degree are oversampled. This is the basic insight upon which the multiplicity adjustment is based. Consequently, to the extent that sampling is affected by differences in degree, respondents must be weighted inversely by their degree. This suggests that the formula for degree component must take the following form, where  $DC^i$  is the degree component for individual  $i$ ,  $D^i$  is the actor's self-reported network size, and  $K$  is a positive constant:

$$(35) \quad DC^i = K \frac{1}{D^i}$$

It is also useful for the sampling weights to sum to the sample size, that is,

$$(36) \quad \sum_{i=1}^n W^i = n$$

These two constraints provide the basis for deriving an expression for the individualized degree components. From equation 29, and employing the individualized version of the degree component (i.e.,  $DC^i$  for respondent  $i$ ), and the individual values for the recruitment component (i.e., where individual  $i$  is a member of group  $X$ ,  $RC^i = RC_X$ ), the above expression expands to

$$(37) \quad \sum_{i=1}^n (DC^i \cdot RC^i) = n$$

Expanding the  $DC$  term from equation 35 yields

$$(38) \quad \sum_{i=1}^n \left( K \frac{1}{D^i} RC^i \right) = n$$

which can be rearranged, providing the means to calculate  $K$ ,

$$(39) \quad K = \frac{n}{\sum_{i=1}^n \left( \frac{1}{D^i} RC^i \right)}$$

The inclusion of this constant, it should be noted, does not affect estimates such as prevalence estimates, but it nonetheless has utility because when frequency distributions are calculated, the number of cases corresponds to the cases for which valid data are available. Substituting equation 39 into equation 35 yields the expression for any respondent *i*'s degree component,

$$(40) \quad DC^i = \frac{n}{\sum_{j=1}^n \frac{1}{D^j} RC^j} \cdot \frac{1}{D^i}$$

Using this expression, the degree component vector is calculated,  $DC = (DC^1, DC^2, DC^3, \dots, DC^n)$ , with values assigned consistent with each respondent's self-reported degree, as well as the constant *K*. In the jazz data set, degree components vary from a minimum of 0.125 for the respondent with the largest network (850) to 5.302 for the respondent with the smallest network (20)—a more than 40-fold variation in values.

Based on the individualized degree component and the recruitment component, the dual-component RDS sampling weight,  $DW^i$ , can be calculated as follows: for individual *i*,

$$(41) \quad DW^i = DC^i \cdot RC^i$$

Expanding the degree component from equation 35 above yields

$$(42) \quad DW^i = \frac{1}{D^i} \cdot K \cdot RC^i$$

By inspection, it is apparent that this dual-component sampling weight combines a multiplicity adjustment (1/*D*) with an adjustment for differential recruitment (*RC*).

The dual-component population estimate is calculated in the standard manner for calculating means for weighted data; that is, for any group *X*, with number of cases  $n_X$ ,

and where  $V^i$  is the variable's value for respondent  $i$  from group  $X$ ,

$$(43) \quad \widehat{DP}_X = \frac{\sum_{i=1}^{n_X} DW^i V^i}{\sum_{i=1}^{n_X} DW^i}$$

This does not represent a distinct RDS population estimator, for when the categorical and individualized weights are employed to estimate the same parameter, the estimates are equal. The difference lies, not in the constituent calculations, but rather in the order in which the calculations are carried out. (For a proof of equivalence, see Appendix C.)

An advantage of the individualized weights is that they permit analysis of continuous variables. Figure 1 displays the analysis of two continuous variables, age and degree, comparing the unweighted estimate, the multiplicity estimate, and a set of dual-component estimates. The cumulative distribution was calculated for each variable. To more clearly display the differences among the estimates, the unweighted estimate was used as a baseline (i.e., each weighted estimate was subtracted from the unweighted estimate). Consequently, the area under the curve does not equal one; instead it indicates the difference between the unweighted estimate and each type of weighted estimate. Consider first the analysis of age (Figure 1A). The weighted estimates consistently exceed the unweighted estimate because of a positive correlation between age and degree ( $r = 0.268$ ), so weighting inflated the estimated number of younger musicians. Multiplicity weighting has substantial effects; for example, the estimated percentage of respondents aged 40 or less is 9% greater in the multiplicity than in the unweighted estimate.

Figure 1A also displays alternative dual-component estimates that differ based on how the continuous variable is partitioned. In contrast with the degree component, the

recruitment component is defined only at the group level, because it is calculated based on aggregate recruitment patterns. Therefore, a step when calculating weights for a continuous variable is to choose an appropriate level of aggregation— that is, whether to partition the sample at the median or by terciles, quartiles, or a finer gradation. The rationale for identifying an appropriate level of aggregation is that any aggregation level omits within-category information regarding differential recruitment, so lower aggregation levels capture less information than higher levels. However, a too-high aggregation level subdivides the sample into so many cells that many have small or zero values, thereby producing computational instability.

Some guidance can be drawn from other statistical procedures for which rules have been devised regarding cell size; for example, for chi square that number is a minimum of five. Employing the chi-square rule, an aggregation level of three is appropriate for both continuous variables in Figure 1, because if divided by quartiles, the minimum cell sizes would be less than 5 (i.e., 4 and 2, respectively, for degree and age), but if divided by terciles, all cells have values in excess of 5 (i.e., 18 and 7, respectively, for degree and age). As defined in this manner, the appropriate aggregation level may vary across continuous variables, with lower levels when recruiter and recruit attributes on the variable are correlated, because cases would cluster on the principal diagonal with smaller numbers of cases elsewhere in the matrix, and higher levels when recruiter and recruit attributes are independent, because cases would be uniformly distributed throughout the recruitment matrix. From the standpoint of RDS analysis, this is a less than ideal procedure, because estimates remain mathematically stable despite low or zero values in cells of the recruitment matrix; what is more significant is the mean cell size—

that is, the number of recruits and recruiters, respectively, in each subgroup. Therefore, a simpler approach—specifying an optimal mean number of cases per cell—appears more appropriate. Where  $AL$  is the aggregation level,  $n$  is the sample size, and  $n_C$  is the mean number of cases per cell, the equation for  $n_C$  is

$$(44) \quad n_C = \frac{n}{AL^2}$$

Solving this equation for aggregation level then yields

$$(45) \quad AL = \sqrt{\frac{n}{n_C}}$$

Figure 1A displays dual-component estimates for cumulative age using a range of aggregation levels from two to nine. Levels two, three, and four deviate to a progressively greater degree from the multiplicity estimate. This suggests that levels two and three are too crude to adequately capture recruitment patterns. Intermediate levels of aggregation from four through six are highly convergent, with correlations among the estimates of 0.998. This range of levels defines what may be termed a *zone of convergence* and corresponds to mean cell sizes from 6.9 for the aggregation level of 6 ( $251/6^2$ ) to 15.7 for the aggregation level of four ( $251/4^2$ ). Higher aggregation levels show signs of instability, as indicated by the nonmonotonic nature of the differences among levels 7 through 9. For example, the estimates for the aggregation level of seven are intermediate between those for levels eight and nine. The intermediate range of levels, those falling within the zone of convergence, are optimal in that they best avoid both instability from too high an aggregation level, and loss of information from a too-low level. As thus defined, the optimal aggregation minimizes, but cannot eliminate, loss of information.



Analyses of continuous variables from several RDS data sets suggest that this data set is not unique. The zone of convergence tends to occur for mean cell sizes in the range  $12 \pm 4$ , with higher levels producing instability and lower levels falling between the convergence zone and the multiplicity estimate. However, additional research will be required to determine the appropriateness of this guideline for a larger range of data sets and continuous variables. What would be especially helpful would be an analytically derived procedure for confirming the presence of a zone of convergence and calculating its boundaries. This will be the topic of a future paper. In the meantime, a test for convergence is useful; analyses should employ a range of aggregation levels to confirm the presence of a convergence zone. Further research would also be useful to devise means for making computations involving higher levels of aggregation numerically stable.

The results of the degree analysis differ from the age analysis in three respects. First, the effect of weighting is greater: both weighted estimates differ by more than 33% from the unweighted estimate. This is to be expected, because both estimates include a multiplicity adjustment, in which respondents of greatest degree are given the smallest weight. A second difference from the age analysis is that dual-component estimates with varying aggregation levels are highly convergent, to such an extent that they could not be clearly distinguished, so only the aggregation level of five is displayed. Third, both the multiplicity and the dual-component estimates are convergent, with a maximum difference of 0.96% for respondents of degree 125. These convergences, among dual-component estimates and between them and the multiplicity estimate, suggest that in this case the effects of differential recruitment on degree estimates are minor.

There are theoretic reasons to suppose that a feature of RDS survey design may weaken the effects of differential recruitment on degree estimation. Recruitment quotas limit the number of peers a respondent can recruit, generally to no more than three. Consequently, even respondents with small networks can fulfill the quota, so the correlation between degree and number of recruits tends to be small—for example, 0.051 in the original RDS study (Heckathorn 1997) and  $-0.044$  in the NYC jazz study. Therefore, an essential element of differential recruitment bias—differential recruitment effectiveness—is lacking. Because respondents of differing degree have equal opportunities to recruit, the lack of a substantial difference between the multiplicity and the dual-component estimates for degree is expected. However, alternative and potentially useful research designs that are discussed in the conclusion could make this bias quite large. Consequently, the ability to control for this source of bias is useful because it expands the range of viable research designs.

### **III. Controlling for Differential Recruitment Bias in Degree Estimation**

The dual-component approach presented in the previous section did not introduce a new population estimator. For the manner in which the degree and recruitment components were defined guaranteed their equivalence to the original formulation. That is, the sample weight,  $W = \hat{P}/C$  was divided into two components,  $\hat{P}/\hat{E}$  and  $\hat{E}/C$  such that, when multiplied, the intermediate term cancels out, i.e.,  $\hat{P}/\hat{E} * \hat{E}/C = \hat{P}/C$ . This equivalence was not altered by the algebraic manipulations that lead to the individualized sampling weights, however different the expression may appear. Consequently, the dual-component estimator can be validly expected to have the same properties as the original RDS estimator, i.e., to be asymptotically unbiased when assumptions 1 to 6 above are

satisfied. In contrast, this section introduces a new estimator intended to control for bias resulting from violation of assumption six, differential recruitment.

The new estimator is based on the ability to calculate degree estimates in a manner that controls for differential recruitment by degree. Using the weights derived from the dual-component analysis of degree, degree estimates can be derived for groups defined by other variables, such as gender or airplay. Using the standard method for calculating means using weighted variables—that is, where  $n_X$  is the number of cases in any group X,  $D^i$  is degree of individual I from group X, and  $DWD^i$  is the dual-component weight for the individual's degree—the adjusted degree estimate for the group,  $\widehat{AD}_X$ , is

$$(46) \quad \widehat{AD}_X = \frac{\sum_{i=1}^{N_X} (D^i \cdot DWD^i)}{\sum_{i=1}^{N_X} (DWD^i)}$$

For example, the degree estimate for females increases from the multiplicity estimate of 102.566 to 103.849 (see Table IA).

The expression for adjusted degree can be simplified to more clearly reveal its relationship to the multiplicity approach to degree estimation. From equation 42 and where  $RCD^i$  is the recruitment component for degree for the degree group into which individual i falls, equation 46 can be expanded as follows:

$$(47) \quad \widehat{AD}_X = \frac{\sum_{i=1}^{n_X} (D^i \frac{1}{D^i} K \cdot RCD^i)}{\sum_{i=1}^{n_X} (\frac{1}{D^i} K \cdot RCD^i)}$$

By algebraic manipulation, this expression can be simplified because the D terms in the numerator cancel one another, and K is a constant that appears in both the numerator and

the denominator, so the reduced expression is

$$(48) \quad \widehat{AD}_X = \frac{\sum_{i=1}^{n_X} RCD^i}{\sum_{i=1}^{n_X} \left(\frac{1}{D^i} RCD^i\right)}$$

This expression can be seen as a weighted version of equation 21's use of a multiplicity adjustment to estimate degree, where the weight is the recruitment component derived from the analysis of degree. It is also apparent that when degree's recruitment component is neutral (i.e.,  $RCD = 1$ ), equations 21 and 48 are equivalent, because in the latter equation the numerator, the sum of recruitment components, becomes equivalent to the number of cases, and in the denominator, the  $(1/D)*RCD$  reduces to  $(1/D) * 1 = 1/D$ . Similarly, if the value of RCD is equal within a group (i.e., for all members a group,  $RCD^i$  has the same value), the effects of RCD cancel so the multiplicity and adjusted estimates are equivalent. This does not occur for groups defined by gender, airplay, or age, because individuals of diverse degree occur within each group. But it is necessarily the case for groups categorized by degree, e.g., in Table IIB's analysis, all respondents in the 20 to 90 degree group have the same RCD, 0.991, and the same is true for the other four groups. As a result, for the degree variable, the multiplicity and adjusted degree estimates are equivalent (see Table IIB). In contrast, because RCD is non-neutral (i.e.,  $RCD \neq 1$ ), the multiplicity and RDS population estimates differ, e.g., the multiplicity estimate of the proportion of the population in the degree 20 to 90 group is 0.539, whereas the RDS estimate is 0.534.

Based on this adjusted degree estimate, a new population estimator can be derived by substituting the adjusted degree estimate for the multiplicity-based estimate in the

original RDS estimator, equation 10; that is, for groups X and Y,

$$(49) \quad \widehat{AP}_X = \frac{\widehat{S}_{YX} \widehat{AD}_Y}{\widehat{S}_{YX} \widehat{AD}_Y + \widehat{S}_{XY} \widehat{AD}_X}$$

When this term is expanded by substituting equation 48's expression for adjusted degree, where  $n_x$  is the number of cases in group X, and  $n_y$  is the number of cases in group Y, the result is

$$(50) \quad \widehat{AP}_X = \frac{\widehat{S}_{YX} \left( \frac{\sum_{j=1}^{n_y} RCD^j}{\sum_{j=1}^{n_y} \left( \frac{1}{D^j} RCD^j \right)} \right)}{\widehat{S}_{YX} \left( \frac{\sum_{j=1}^{n_y} RCD^j}{\sum_{j=1}^{n_y} \left( \frac{1}{D^j} RCD^j \right)} \right) + \widehat{S}_{XY} \left( \frac{\sum_{i=1}^{n_x} RCD^i}{\sum_{i=1}^{n_x} \left( \frac{1}{D^i} RCD^i \right)} \right)}$$

What is notable about this expression is that the recruitment component for the degree variable enters into the degree estimation process for analysis of all other variables. In this way, the estimator compensates for differential recruitment by degree. The adjusted population estimates are displayed in the bottom panels of Tables I and II. The estimate for the proportion of females changes from 0.2381 to 0.2380, an adjustment that changes the estimate only at the fourth decimal place. The effects of adjustment on the estimates for airplay are greater, consistent with the greater dependence of this estimate on degree differences between groups (see Table IB). The estimated proportion with airplay changes from 0.752 to 0.751. The substantial changes in degree estimates and the comparatively small changes in population estimates occur because the population estimates depend not on absolute but on relative degrees; that is, an order-preserving linear transform of groups' estimated degrees has no effect on the population estimate, so

only relative degrees estimates are altered, and these alterations are minor.

Finally, the adjusted population estimate provides the basis for adjusting the sampling weight; that is for any group X, the adjusted sampling weight,  $AW_x$ , is

$$(51) \quad AW_x = \frac{\widehat{AP}_x}{C_x}$$

Like the standard RDS sampling weight, this is a group-level weight. The adjustment process can be extended to the dual-component weight, which is individualized based on each respondent's degree by multiplying the latter by the ratio of the adjusted and the standard weight. That is, for individual  $i$  from group X, the adjusted dual-component weight,  $ADW^i$ , is

$$(52) \quad ADW^i = DW^i \frac{AW_x}{W_x} = DW^i \frac{\left( \frac{\widehat{AP}_x}{C_x} \right)}{\left( \frac{\widehat{P}_x}{C_x} \right)}$$

Equivalently, given that the C terms cancel, a simpler expression for the adjusted dual-component weight for individual  $i$  from group X is

$$(53) \quad ADW^i = DW^i \frac{\widehat{AP}_x}{\widehat{P}_x}$$

This expression permits the adjustment process to be extended to continuous variables. For example, when imported into a program such as STATA, mean values and other statistics can be calculated for continuous variables.

TABLE III ABOUT HERE

To assess the variation of the degree variable's recruitment component, it is useful to examine the magnitude of this term in multiple RDS data sets. Magnitude is defined as

the mean of the absolute differences from unity of the degree recruitment component, i.e., for a degree variable partitioned into  $M$  categories where  $RCD_i$  is the recruitment component for category  $i$ , the magnitude is,  $\sum_{i=1}^M |RCD_i - 1| / M$ ; for example, for the two categories 1.3 and 0.9, the magnitude is  $(|1.3 - 1| + |0.9 - 1|) / 2 = (0.3 + 0.1) / 2 = 0.2$ . Table III shows this term for eight RDS data sets. The magnitude has a minimum value of 0.012 in a study of drug users in New York City (Abdul-Quader et al. 2006), a value of 0.028 in the New York City jazz musician study, whose data are employed in this article, and a maximum of 0.097 in a study of Middletown, Connecticut, injection drug users (Heckathorn 1997).

FIGURE 2 ABOUT HERE

To assess the sensitivity of RDS estimates to change in the magnitude of RCD, analyses were carried out to explore the effects of varying that magnitude. Using the jazz musician data set, the magnitude of RCD was multiplied by integers from zero to 10. This yielded 11 data sets with magnitudes varying from zero (when the multiplier is zero) to 0.028 (when the multiplier is one—the value in the original data set) to a maximum of 0.28 (when the multiplier is 10). Figure 2 shows the effects of these alterations in the magnitude of RCD for two variables, gender and airplay, where the vertical axis is the difference between the original RDS and the adjusted estimates. By inspection, it is apparent that the relationship is approximately linear, with a steeper slope for airplay, the variable whose weight principally depends on its degree component, and a more gentle slope for gender, the variable whose weight is principally determined by its recruitment component. When the multiplier is zero (i.e., the left of Figure 1), there is no differential recruitment by degree, so the standard RDS and adjusted estimates are equal. When the

multiplier is one and the magnitude of RCD is 0.028, the differences between the estimates are those reflected in Tables IA and IB—that is, 0.0001 for gender and 0.0009 for airplay. When the multiplier is 10, the estimates increase 10-fold, to 0.001 for gender and .009 for airplay. Consequently, increasing the magnitude of RCD by an order of magnitude, to an amount nearly triple that found in any of the listed RDS data sets, has rather modest effects.

However, changes in research design that induce an association between degree and opportunities to recruit would produce much larger potential effects. These effects were explored by transforming the NYC jazz data set in a manner consistent with a design intended to increase representation of an otherwise under-sampled group, respondents of small degree. Specifically, recruitment for the quintile of smallest degree, those of degree 20 to 90, was tripled, thereby multiplying by three each entry in the top row of Table IIB's recruitment matrix, with no other changes in data structure, e.g., the new hypothetical recruits were assumed to have the same degree as the actual recruits. The result is a substantial increase in the magnitude of the recruitment component for degree, from 0.028 to 0.246. Based on Figure 2's analysis of the effects of a uniform alteration in this magnitude, the expected bias would be approximately 0.8% for a variable such as airplay for which the degree component is most significant, and 0.08% for a variable such as gender for which the degree component is less significant. However, when standard RDS estimates are generated, the discrepancies are quite different, the gender estimate changes by 4%, whereas the airplay estimates change by the expected 1%. This apparent anomaly results from the non-uniform alteration in the recruitment component by degree. For the lowest quintile, the recruitment component by



degree changes from a near neutral 0.991 to 0.582, for a decrease of 0.409. In contrast, for the next quintile, respondents of degree 100 to 125, it changes from 0.937 to 1.164, for an increase of 0.227. The effects on the gender variable are greater, because though both genders are similar in mean degree, the variance in degree for males is greater so more males lie in the extremes of the degree distribution, including the bottom quintile. Of course, the adjusted estimates remain equivalent in the original and the altered data sets, because the adjustment procedure filters out the confounding effects of differential recruitment by degree. As this hypothetical example demonstrates, the effect of differential recruitment by degree can be both substantial and complex when recruitment effectiveness is associated with degree.

### **Conclusion**

Sampling weights in RDS differ from standard weights because each respondent's weight varies as analyses shift from variable to variable. Employing standard associational methods to derive a uniform set of sampling weights that could be used for all variables or sets of variables would simplify the analyses. However, such an approach would involve a loss of precision. For given that the RDS estimator is asymptotically unbiased for data sets that fit the method's assumptions (Salganik and Heckathorn, 2004), and that those estimates yield different weights for different variables, any method that produced uniform weights would introduce a bias.

The potential for this variable-dependent bias derives from interactions between the two factors that determine weights. These can multiply one another, as when a group of larger degree has a bias toward in-group ties and recruits more effectively; here differentials in degree and differential recruitment both combine to increase over-

sampling of the group. Alternatively, the two factors may counter one another, as when a group of smaller degree is recruited more frequently by the more effectively recruiting groups. In that case, differentials in degree and differential recruitment can either cancel one another to produce a self-weighting sample, or one or the other factor may prove stronger so weights must compensate for either under or over-sampling, respectively. Weighting that is affected by differential recruitment necessarily varies across variables, because differential recruitment is based on network properties.

Socially irrelevant variables, such as having been born in an odd or an even month, do not affect affiliation, and consequently fit a network structure consistent with random mixing. In that case, differential recruitment would be minimal, deriving only from stochastic variation. In contrast, variables that reflect network segmentation, such as race/ethnicity and other demographic factors, can produce patterns of differential recruitment that vary in complex ways across variables. Consequently, the ability of an estimator to control for differential recruitment effects is broadly relevant.

In contrast, controlling for differential recruitment in degree estimation may have only minor effects when standard RDS research protocols are followed because recruitment quotas induce a small correlation between degree and recruitment effectiveness. Nevertheless, a conservative research approach requires that this expectation be confirmed.

Moreover, the ability to control for differential recruitment in degree estimation increases the range of research protocols from which unbiased results can be expected. For example, changes in recruitment quotas could enable the sampling process to adapt to less dense networks. Network density is seldom an issue when sampling drug users or

gay men because both populations tend to be gregarious. In contrast, the network structure of commercial sex workers is more diverse. RDS has been successfully used with sex workers in two cities in Vietnam (Johnston et al. 2006). However, it performed poorly in two cities in Estonia and Russia (Simic et al. 2006), for reasons that may include limitations in the ability of prostitutes to independently form networks. When sampling is without replacement (i.e., respondents can be interviewed only once), networks are sparse, and when a respondent has made the maximum number of allowable recruitments, others to whom the respondent is connected may become stranded, with no network connections linking them directly or indirectly to any potential recruiters. In essence, the recruitment process breaks what had initially been a large low-density component into multiple isolated components. This can be avoided by a change in research design, in which recruitment would take place in two stages. The first would employ a uniform modest quota, as is standard practice, until the sample had attained adequate sociometric depth (i.e., number of waves) to ensure that the population's network had been adequately penetrated. Second, those respondents who had fulfilled their quotas would then be given an equal number of additional recruitment rights, and this process would be repeated a specified number of times until the target sample size was reached. In this way, respondents located in isolated components created during the first stage could potentially be reached during the second stage, thereby increasing the sampling method's performance in low-density networks. An effect of this alteration in research design would be to increase the association between degree and recruitment effectiveness, thereby potentially introducing a large bias from differential recruitment by degree. In that case, having the means to post-stratify the sample to compensate for that

bias may prove important.

A second context in which means for controlling for differential recruitment by degree may become important occurs when the sample is stratified to increase recruitment of groups of special interest. Generally, this involves a larger recruitment quota for members of these groups. The effect is to increase the recruitment effectiveness of the targeted groups. Consequently, if the groups' degree differs from the norm, differential recruitment by degree will occur.

The relevance of the weighting procedures introduced in this paper differs based on the form of analysis. When point estimates are at issue, weighting is always potentially important, as illustrated by an RDS study of HIV prevalence in San Francisco and Chicago (Ramirez-Valles et al. 2005a). In San Francisco, recruitment by HIV status was nearly identical for both HIV positives and negatives. Similarly, degree was unrelated to HIV status, so the sample was self-weighting. That is, the sample composition (46.1%) approximated the RDS estimate (48.7%). In contrast, in Chicago the differences in recruitment patterns were substantial, as were differences in recruitment effectiveness, so both elements required for differential recruitment effects were present. Moreover, HIV positives had networks nearly twice as large as negatives, so weights had a very substantial effect. Prevalence in the sample was 24.7%, but the RDS estimate was 16.8%. Point estimates for continuous variables also require weights. For example, the mean age for New York City female jazz musicians varies from 47.9 without weights to 44.7 for the multiplicity weight and 43.5 for the dual-component weight.

Weights function differently in multivariate analysis. It is now recognized

(Winship and Radbill 1994) that weighting frequently has little effect on regression analyses, because these depend on correlations among variables that tend to change only slightly when weights inflate or deflate the value of a variable. Consequently, Winship and Radbill recommend conducting the analysis with weights, and then repeating the analysis with no weights. If the results of the two analyses are convergent, they recommend reporting the unweighted result. The advantage of this procedure is that weighting produces wider confidence intervals, so they should be employed only when necessary. Of course, that determination requires the ability to replicate analyses both with and without weights, for which the sort of weighting procedure introduced in this paper is needed. (For an example of this procedure that uses RDS data, see Ramirez-Valles, et al. in press.)

Further development of RDS in several directions would be useful. First, the effects of differential recruitment on variance estimation should be explored (for a detailed treatment of bootstrap confidence intervals in RDS analysis, see Salganik 2006). When some groups recruit more effectively than others, and groups vary in the variance in their recruitment behaviors, an approach to variance estimation that assumes uniform recruitment effectiveness would fail to include that source of variance. Second, whether efforts to derive variance estimates analytically (Volz and Heckathorn, in press) will prove compatible with the dual-component approach is a more difficult question that requires additional research. It also remains to be seen whether statistical packages, such as SUDAAN, that are designed to accommodate highly complex weighting systems can be adapted to the dual-component approach. In any case, the introduction of a new estimation procedure requires corresponding adjustments in procedures for variance

estimation.

---

**Table IA. RDS and Multiplicity Estimates for Gender, NYC Jazz Musicians**

---

Gender of person who recruited	Gender of recruit		Total recruits by each group (RB)
	Male	Female	
Male			
Recruitment count	127	25	152
(Recruitment proportions)	(0.836)	(0.164)	1
Female			
Recruitment count	51	40	91
(Recruitment proportions)	(0.56)	(0.44)	1
Total recruits of each group (RO)	178	65	243
Sample composition, (including seeds)	0.737	0.263	1
Mean degree, (multiplicity estimate)	109.225	102.566	
Population estimate, (multiplicity estimate)	0.7267	0.2752	1
Equilibrium proportion	0.773	0.227	1
Sampling weight	1.033	0.907	
Population estimate, (standard RDS estimate)	0.7619	0.2381	1
Degree component	0.985	1.049	
Recruitment component	1.048	0.864	
Mean degree, (adjusted estimate)	110.513	103.849	
Population estimate, (adjusted estimate)	0.7620	0.2380	1

---

**Table IB. RDS and Multiplicity Estimates for Airplay, NYC Jazz Musicians**

Airplay of person who recruited	Airplay of recruit		Total by each group (RB)
	Yes	No	
Yes			
Recruitment count	155	33	188
(Recruitment proportion)	(0.834)	(0.176)	1
No			
Recruitment count	40	11	51
(Recruitment proportion)	(0.784)	(0.216)	1
Total recruits of each group (RO)	195	44	239
Sample composition, (including seeds)	0.822	0.178	1
Mean degree, (multiplicity estimate)	116.66	79.074	
Population estimate, (multiplicity estimate)	0.759	0.241	1
Equilibrium proportion	0.817	0.183	1
Sampling weight	0.914	1.396	
Population estimate, (standard RDS estimate)	0.752	0.248	1
Degree component	.92	1.357	
Recruitment component	0.994	1.028	
Mean degree, (adjusted estimate)	118.174	79.717	
Population estimate, (adjusted estimate)	0.751	0.249	1



**Table IIA. RDS and Multiplicity Analysis for Age, Partitioned by Quintiles**

Age of person who recruited	Age of recruit					Total
	Recruitment count (recruitment proportion)					
	Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101	
Age 20-33	13 (0.433)	10 (0.333)	3 (0.1)	2 (0.067)	2 (0.067)	30 1
Age 34-42	15 (0.268)	17 (0.304)	12 (0.214)	9 (0.161)	3 (0.054)	56 1
Age 43-49	7 (0.14)	9 (0.18)	8 (0.16)	14 (0.28)	12 (0.24)	50 1
Age 50-58	7 (0.117)	9 (0.15)	17 (0.283)	13 (0.217)	14 (0.233)	60 1
Age 59-101	8 (0.148)	4 (0.074)	10 (0.185)	15 (0.278)	17 (0.315)	54 1
Total recruits of each group (RO)	50	49	50	53	48	250
Sample composition, (including seeds)	0.202	0.198	0.198	0.205	0.198	1
Mean degree, (multiplicity estimate)	82.144	108.176	109.209	97.598	183.059	
Population estimate, (multiplicity estimate)	0.260	0.193	0.192	0.222	0.114	1
Equilibrium proportion	0.233	0.219	0.186	0.192	0.17	1
Sampling weight	1.49	1.083	0.91	1.011	0.497	
Population estimate, (standard RDS estimate)	0.3	0.214	0.18	0.207	0.098	1
Degree component	1.287	0.977	0.968	1.083	0.577	
Recruitment component	1.158	1.108	0.94	0.933	0.861	
Mean degree, (adjusted estimate)	82.81	109.777	110.549	98.569	186.183	
Population estimate, (adjusted estimate)	0.301	0.213	0.18	0.208	0.098	1

**Table IIB. RDS and Multiplicity Analysis of Degree, Partitioned by Quintiles**

Degree of person who recruited	Degree of recruit					
	Recruitment count (recruitment proportion)					
	Degree 20-90	Degree 100-125	Degree 150-200	Degree 220-400	Degree 450-850	Total
Degree 20-90	8 (0.296)	4 (0.148)	8 (0.296)	4 (0.148)	3 (0.111)	27 1
Degree 100-125	7 (0.163)	9 (0.209)	9 (0.209)	9 (0.209)	9 (0.209)	43 1
Degree 150-200	16 (0.267)	15 (0.25)	13 (0.217)	10 (0.167)	6 (0.1)	60 1
Degree 220-400	5 (0.104)	5 (0.104)	15 (0.313)	13 (0.271)	10 (0.208)	48 1
Degree 450-850	3 (0.115)	6 (0.231)	6 (0.231)	7 (0.269)	4 (0.154)	26 1
Total recruits of each group (RO)	39	39	51	43	32	204
Sample composition, (including seeds)	0.198	0.202	0.247	0.202	0.152	1
Mean degree, (multiplicity estimate)	40.019	101.658	175.758	308.057	539.698	
Population estimate, (multiplicity estimate)	0.539	0.217	0.154	0.072	0.031	1
Equilibrium proportion	0.196	0.189	0.253	0.209	0.154	1
Sampling weight	2.705	1.008	0.637	0.367	0.204	
Population estimate, (standard RDS estimate)	0.534	0.203	0.157	0.074	0.031	1
Degree component	2.731	1.075	0.622	0.355	0.202	
Recruitment component	0.991	0.937	1.025	1.034	1.009	
Mean degree, (adjusted estimate)*	40.019	101.658	175.758	308.057	539.698	
Population estimate, (adjusted estimate)*	0.534	0.203	0.157	0.074	0.031	1

\*Note: The multiplicity and adjusted degree estimates are equal, as are the standard and adjusted RDS population estimates, because there is no within-group variation in the value of the degree recruitment component.

---

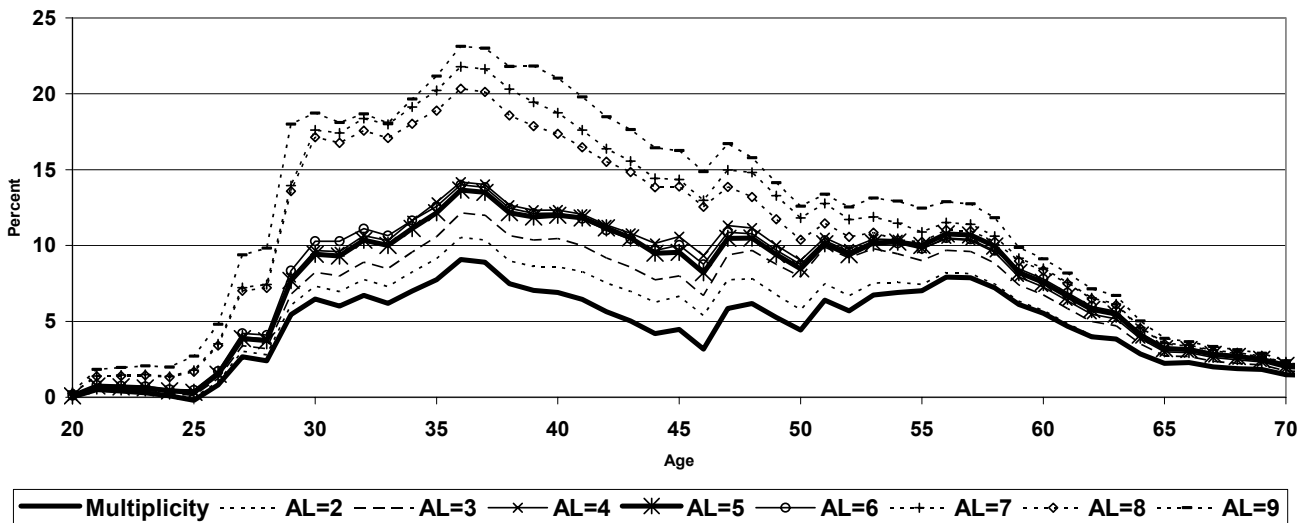
***Table III. Magnitude of the Recruitment Component of Degree (RCD) in 8 RDS Data Sets***

---

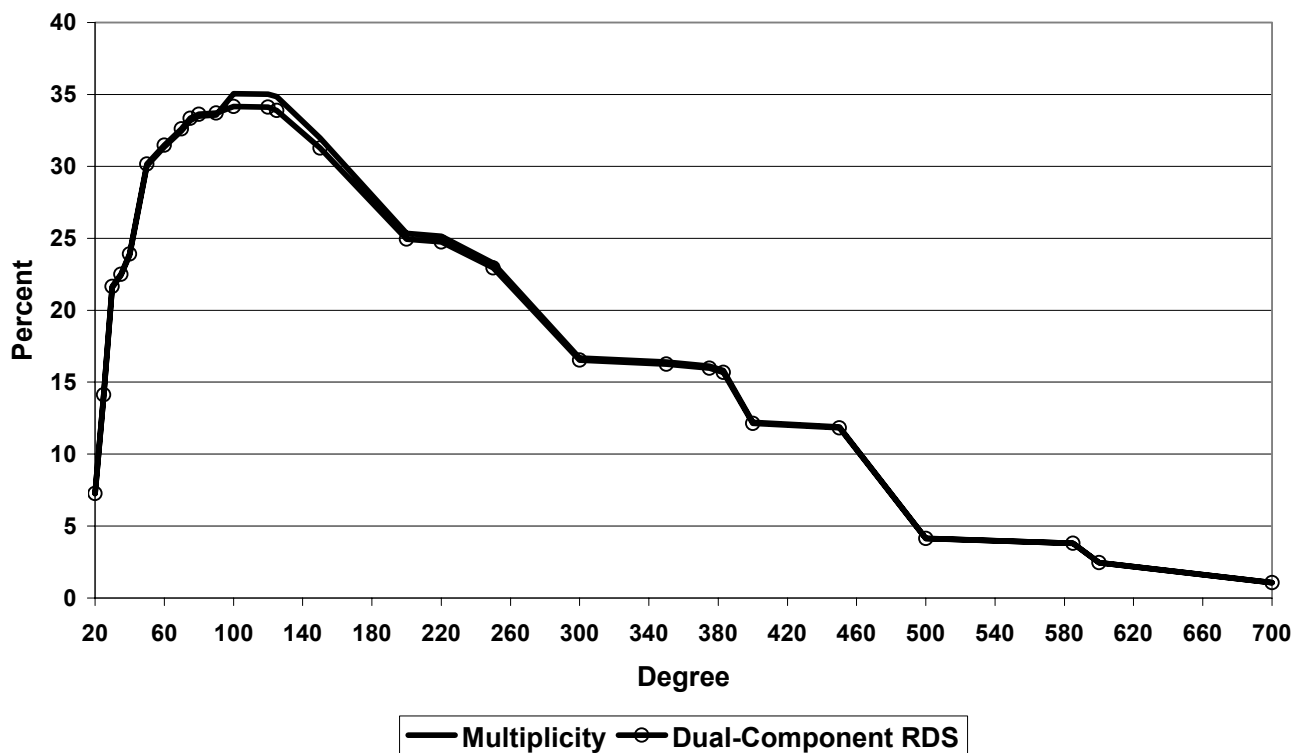
New York City drug users (Abdul-Quader et al. 2006)	0.012
New York City jazz musicians (Heckathorn and Jeffri 2003)	0.028
Chicago Latino gay, bisexual, and transsexual (Ramirez-Valles et al. in press)	0.039
Cornell University Undergraduates (Wejnert and Heckathorn 2005)	0.051
San Francisco Latino gay, bisexual, and transsexual (Ramirez-Valles et al. in press)	0.078
New London (CT) injection drug users (Heckathorn et al. 1999)	0.078
San Francisco jazz musicians (Heckathorn and Jeffri 2003)	0.094
Middletown (CT) injection drug users (Heckathorn 1997)	0.097

---

**Figure 1A: Age -- Multiplicity and Weighted Estimates by Aggregation Level (AL)**  
(Using Unweighted Estimate as a Baseline)

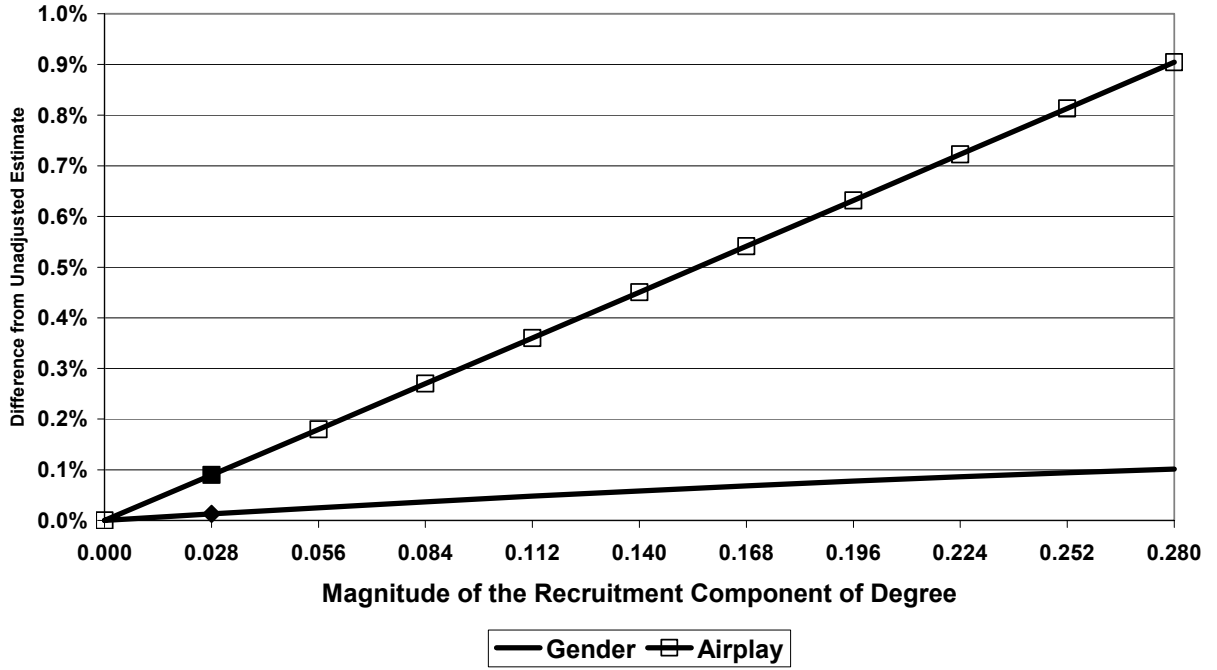


**Figure 1B: Degree -- Multiplicity and Dual-Component RDS Estimate**  
(Using Unweighted Estimate as a Baseline)



**Figure 2: Effects of Varying the Magnitude of Degree's Recruitment Component**

Magnitude ( (RCD = 2.8%) is varied from zero to a multiplier of ten



## REFERENCES

- Abdul-Quader, Abu S., Douglas D. Heckathorn, Courtney McKnight, Heidi Bramson, Chris Nemeth, Keith Sabin, Kathleen Gallagher, and Don C. Des Jarlais. 2006. "Effectiveness of Respondent Driven Sampling for Recruiting Drug Users in New York City: Findings From a Pilot Study." *Journal of Urban Health*, 83:459-476.
- Bonacich, Phillip F. 1972. "Factoring and Weighting Approaches to Status Scores and Clique Identification." *Journal of Mathematical Sociology*. 2:113-120.
- Brewer, K. R. W., and M. Hanif. 1983. *Sampling with Unequal Probability*. New York: Springer-Verlag.
- Brown, Lawrence D., Morris L. Eaton, David A. Freedman, Stephen P. Klein, Richard A. Olshen, Kenneth W. Wachter, Martin T. Wells and Donald Ylvisaker. 1999. "Statistical Controversies in Census 2000." Technical Report 537, Department of Statistics, University of California–Berkeley.
- Cochran, W. G. 1977. *Sampling Techniques*. 3d ed. New York: Wiley.
- Erickson, Bonnie H. 1979. "Some Problems of Inference from Chain Data." *Sociological Methodology* 10:276–302.
- Frank, Ove. 1979. "Estimation of Population Totals by Use of Snowball Sampling." In *Perspectives on Social Network Research*, edited by P.W. Holland and S. Leinhardt. New York: Academic Press.
- Frank, Ove and Tom Snijders. 1994. "Estimating the Size of Hidden Populations Using Snowball Sampling." *Journal of Official Statistics* 10:53–67.
- Goodman, Leo A. 1961. "Snowball Sampling." *Annals of Mathematical Statistics* 32:148–70.

- Hansen, M. H., and W. N. Hurwitz. 1943. "On the Theory of Sampling from Finite Populations." *Annals of Mathematical Statistics* 14(4):333–62.
- Heckathorn, Douglas D. 1997. "Respondent Driven Sampling: A New Approach to the Study of Hidden Populations." *Social Problems* 44:174–99.
- . 2002. "Respondent Driven Sampling II: Deriving Statistically Valid Population Estimates from Chain-Referral Samples of Hidden Populations." *Social Problems* 39: 11-34.
- Heckathorn, Douglas D. and Joan Jeffri. 2003. "Social Networks of Jazz Musicians." Pp. 48-61 in *Changing the Beat: A Study of the Worklife of Jazz Musicians*. Volume III: Respondent-Driven Sampling: Survey Results by the Research Center for Arts and Culture, National Endowment for the Arts Research Division Report #43. Washington, DC.
- Heckathorn, Douglas D., Robert S. Broadhead, Denise L. Anthony and David L. Weakliem. 1999. "AIDS. and Social Networks: Prevention through Network Mobilization." *Sociological Focus* 32:159–79.
- Heckathorn, Douglas D., Salaam Semaan, Robert S. Broadhead, and James J. Hughes. 2002. "Extensions of Respondent-Driven Sampling: A New Approach to the Study of Injection Drug Users Aged 18–25." *AIDS and Behavior* 13(1):55-67.
- Johnston, Lisa Grazina, Keith Sabin, Mai Thu Hien and Pham Thi Huong. 2006. "Assessment of Respondent Driven Sampling for Recruiting Female Sex Workers in Two Vietnamese Cities: Reaching the Unseen Sex Worker." *Journal of Urban Health*. In press.
- Kalton, Graham. 1983. *Introduction to Survey Sampling*. Newbury Park, CA: Sage

- Publications.
- Kemeny, John G. and J. Laurie Snell. 1960. *Finite Markov Chains*. Princeton, NJ: Van Nostrand.
- Klov Dahl, Alden. S. 1989 “Urban Social Networks: Some Methodological Problems and Possibilities.” In *The Small World*, ed. M. Kochen (ed.). New Jersey: Norwood.
- MacKellar, D., L. Valleroy, J. Karon, G. Lemp and R. Janssen. 1996. “The Young Men’s Survey: Methods for Estimating HIV Seroprevalence and Risk Factors among Young Men Who Have Sex with Men.” *Public Health Reports* 111(Supplement):138–44.
- Marsden, Peter V. 1990. “Network Data and Measurement.” *Annual Review of Sociology* 16:435–63.
- Ramirez-Valles, Jesus., Douglas D. Heckathorn, Raquel Vázquez, Rafael M. Diaz and Richard T. Campbell. 2005a. “From Networks to Populations: The Development and Application of Respondent-Driven Sampling among IDUs and Latino Gay Men.” *AIDS and Behavior* 9(4):387–402.
- . 2005b. “The Fit between Theory and Data in Respondent-Driven Sampling: Response to Heimer.” *AIDS and Behavior* 9(4):409–414.
- Ramirez-Valles, Jesus., Dalia Garcia, Richard T. Campbell, Rafael M. Diaz, and Douglas D. Heckathorn. In press. “HIV Infection, Sexual Risk, and Substance Use among Latino Gay and Bisexual Men and Transgender Persons.” *American Journal of Public Health*.
- Rothbart, George S., Michelle Fine and Seymour Sudman. 1982. “On Finding and Interviewing the Needles in the Haystack: The Use of Multiplicity Sampling.”



- Public Opinion Quarterly* 46:408–421.
- Salganik, Matthew J. 2006. “Variance estimation, design effects, and sample size calculations for respondent-driven sampling.” *Journal of Urban Health*, in press.
- Salganik, Matthew J. and Douglas D. Heckathorn. 2004. “Sampling and Estimation in Hidden Populations Using Respondent-Driven Sampling.” *Sociological Methodology* 35:193–238.
- Simic, Milena, Lisa Johnston, Lucy Platt, Sladjana Baros, Violeta Andjelkovic and Tim Rhodes. 2006. “Exploring Barriers to Respondent Driven Sampling in Sex Worker and Drug-Injecting Sex Worker Populations in Eastern Europe.” *Journal of Urban Health*. In press.
- Sirken, M.G. 1970. “Household Surveys with Multiplicity.” *Journal of the American Statistical Association* 65:257–66.
- Snijders, T.A.B. 1992. “Estimation on the Basis of Snowball Samples: How to Weight.” *Bulletin de Methodologie Sociologique* 36:59–70.
- Spreen, Marius. 1992. “Rare Populations, Hidden Populations, and Link-Tracing Designs: What and Why?” *Bulletin de Methodologie Sociologique* 36:34–58.
- Sudman, Seymour and Graham Kalton. 1986. “New Developments in the Sampling of Special Populations.” *Annual Review of Sociology* 12:401–29.
- Thompson, S.K. and O. Frank. 2000. Model-Based Estimation with Linktracing Sampling Designs. *Survey Methodology* 26(1):87–98.
- Volz, Erik and Douglas D. Heckathorn. In press. “Probability-Based Estimation Theory for Respondent-Driven Sampling.” *Journal of Official Statistics*.
- Wang, Jichuan, Robert G. Carlson, Russel S. Falck, Harvey A. Siegal, Ahmmed Rahman

- and Linna Li. 2005. "Respondent-Driven Sampling to Recruit MDMA Users: A Methodological Assessment." *Drug and Alcohol Dependence* 78:147–57.
- Watters, John K. and Patrick Biernacki. 1989. "Targeted Sampling: Options for the Study of Hidden Populations." *Social Problems* 36(4):416–30.
- Watts, Duncan J. 2003. *Small Worlds: The Dynamics of Networks between Order and Randomness*. Princeton, NJ: Princeton University Press.
- Wejnert, Cyprian and Douglas D. Heckathorn. 2005. "Respondent-Driven Sampling and Social Networks: A New Sampling Method." Paper presented at the annual meetings of the American Sociological Association meetings, Philadelphia, PA, August 13–16.
- Winship, Christopher and L. Radbill. 1994. "Sampling Weights and Regression Analysis." *Sociological Methods and Research* 23(2):230–57.

## **Footnotes**

1. For an alternative derivation of an RDS estimator, see Volz and Heckathorn (in press). This approach offers improved analytical tractability and analytical variance estimation, and it provides means for reducing the variance of estimates. It also allows for the estimation of continuous variables; however, the ability to control for differential recruitment bias is limited to nominal variables, a limitation that may be overcome in subsequent work.

## Appendix A. Calculation Procedures for RDS Sampling Weights

The following is a confabulated data set with 20 cases. It consists of the respondent identification (RID) for each respondent and the RID of each respondent's recruiter, along with each respondent's degree (i.e., self-reported network size) and value for a dichotomous variable. Note that respondent 1 lacks a recruiter; this is the seed. Note also that the degree data for respondents 3 and 12 are missing.

Table A.I. Sample Data

RID	RID of recruiter	Degree	Variable
1	NA	8	A
2	1	8	A
3	1	NA	A
4	1	10	B
5	2	5	A
6	2	7	B
7	3	4	B
8	3	7	B
9	3	5	A
10	4	2	B
11	5	4	A
12	5	NA	B
13	5	3	A
14	7	2	B
15	7	3	B
16	7	3	A
17	8	7	B
18	9	3	B
19	9	5	A
20	9	8	B

From these data, the recruitment matrix is constructed by matching recruiters and recruits. For example, respondent 1, a member of group A, recruited respondents 2 and 3, both members of group A, and respondent 4, a member of group B, thereby adding two recruitments from A to A, and one from A to B. The transition probabilities are then calculated based on each group's pattern of recruitment.

Table A.II. Recruitment Matrix			
	Group A recruits	Group B recruits	Tota 1
Person who recruited			1
Group A recruitment count	7	7	14
Group A selection proportion	0.5	0.5	1
Group B recruitment count	1	4	5
Group B selection proportion	0.2	0.8	1
Total Recruitment of Respondents, RO	8	11	19
Recruitment proportion, R	0.421	0.579	1
Sample proportional composition, C (including seeds)	0.45	0.55	1

Note that the sum of cases in the recruitment table is the number of respondents minus the number of seeds; that is,  $20 - 1 = 19$ . This table also displays the recruitment proportions—that is, the proportions derived from the recruitment table, and also the sample proportion, which includes the seeds.

### The Degree Estimate

Degree estimation in traditional RDS analysis is based on a multiplicity adjustment. That is, because respondents are sampled in proportion to their degree, their degrees must be weighted by the inverse. The estimated degree for group A is as follows:

$$(A1) \quad \widehat{D}_A = \frac{n_A}{\sum_{i=1}^{n_A} \frac{1}{D^i}} = \frac{7}{\frac{1}{8} + \frac{1}{5} + \frac{1}{5} + \frac{1}{4} + \frac{1}{3} + \frac{1}{3} + \frac{1}{5}} = 4.264$$

Note that this calculation excludes respondent 1, the seed, because it was not selected through peer recruitment (see Salganik and Heckathorn 2004:215). It also excludes another member of group A, respondent 3, for which degree data are missing. Consequently, though 9 respondents fall within group A, the numerator is 7. Using the

same procedure, the estimated degree for group B is 3.891.

### The Population Estimate

Based on the selection proportions and the degree estimates, the population estimate can be calculated as follows, where  $\hat{D}_x$  is the estimated degree of group X, and  $\hat{S}_{xy}$  is the transition probability from group X to Y:

$$(A2) \quad \hat{P}_A = \frac{\hat{S}_{BA}\hat{D}_B}{\hat{S}_{BA}\hat{D}_B + \hat{S}_{AB}\hat{D}_A} = \frac{0.2 \cdot 3.891}{0.2 \cdot 3.891 + 0.5 \cdot 4.264} = 0.267$$

The population estimate for group B is therefore  $\hat{P}_B = 1 - 0.267 = 0.733$ .

### Sampling Weights

The sampling weight is defined as the ratio of the population estimate, P, and the sample proportional composition, C. Therefore, the sampling weights for each group are calculated as follows:

$$(A3) \quad W_A = \frac{\hat{P}_A}{C_A} = \frac{0.267}{0.45} = 0.594$$

$$W_B = \frac{\hat{P}_B}{C_B} = \frac{0.733}{0.55} = 1.332$$

A vector of these weights can be imported into statistical analysis programs such as STATA or SPSS. The weights for the current data set appear in Table A.III.

Table A.III. Sampling Weights	
RID	Weight
1	0.594
2	0.594
3	0.594
4	1.332
5	0.594
6	1.332
7	1.332

8	1.332
9	0.594
10	1.332
11	0.594
12	1.332
13	0.594
14	1.332
15	1.332
16	0.594
17	1.332
18	1.332
19	0.594
20	1.332

For example, the following is the output from SPSS comparing an unweighted and a weighted analysis. The former reflects merely the sample composition; the latter coincides with the above-derived RDS population estimates. Note also that whereas in the recruitment matrix the sum of cases is 19, in Table A.IV the sample size is 20. This is based on a decision to assign to the seed the sampling weight appropriate to its category, group A.

Table A.IV. Unweighted and Weighted Analyses

<b>Unweighted</b>					
		Frequency	Percentage	Valid percentage	Cumulative percentage
Valid	A	9	45	45	45
	B	11	55	55	100
	Total	20	100	100	
<b>Weighted</b>					
		Frequency	Percentage	Valid percentage	Cumulative percentage
Valid	A	5.3	26.7	26.7	26.7
	B	14.7	73.3	73.3	100
	Total	20	100	100	

A more conservative approach, though one that would entail loss of data, would be to calculate the weight using the recruitment proportion rather than the sample proportion (i.e., group A's weight would then be  $0.267/0.421 = 0.634$ ), and assign a

weight of zero to the seeds; in that case the sample size would be equivalent to that in the recruitment table. The decision to use the former approach is based on the judgment that data loss should be accepted only when inclusion of problematic data would introduce more than trivial amounts of bias, a possibility that appears highly unlikely in this context, where seeds generally compose only a modest proportion of a typical RDS sample.

### **Dual-Component Sampling Weights**

Dual-component weights are calculated based on four terms: the population estimate, the sample composition, the equilibrium for each group, and each respondent's degree.

#### ***Recruitment Component***

The recruitment component of the dual weight is calculated from the sample composition,  $C$ , and the Markov equilibrium,  $E$ . The equilibrium calculated from the transition probabilities for group A is

$$(A4) \quad \widehat{E}_A = \frac{\widehat{S}_{BA}}{\widehat{S}_{BA} + \widehat{S}_{AB}} = \frac{0.2}{0.2 + 0.5} = 0.286$$

Similarly, the equilibrium for group B is  $1 - 0.286 = .714$ .

The recruitment components for groups A and B are therefore

$$(A5) \quad RC_A = \frac{\widehat{E}_A}{C_A} = \frac{0.286}{0.45} = 0.635$$

$$RC_B = \frac{\widehat{E}_B}{C_B} = \frac{0.714}{0.55} = 1.299$$

#### ***Degree Component***

The group-based degree component is the ratio of the group's population estimate



and equilibrium; that is,

$$(A6) \quad DC_A = \frac{\widehat{P}_A}{\widehat{E}_A} = \frac{0.267}{0.286} = 0.936$$

$$DC_B = \frac{\widehat{P}_B}{\widehat{E}_B} = \frac{0.733}{0.714} = 1.026$$

However, what is useful for tasks such as analyzing continuous variables is the individualized version of the degree component (equation 40 in the text). The degree component for a respondent is the product of two terms, the reciprocal of the respondent's degree (i.e., a multiplicity adjustment) and a constant K, which is calculated from the recruitment components and degrees.

This constant is calculated from three terms: the recruitment component, each respondent's degree, and the number of cases for which degree information is available. That is, where n is the number of valid cases for the variable being analyzed, where for any individual i,  $RC^i$  is the recruitment component assigned to the individual's group (i.e., for individual i from any group X,  $RC^i = RC_X$ ), and  $D^i$  is the individual's degree, the constant K is

$$(A7) \quad K = \frac{n}{\sum_{i=1}^n \frac{1}{D^i} RC^i}$$

This constant is an estimate for the overall mean degree for respondents in the system.

The calculation of K for the sample data is shown in Table A.V.

Table A.V. Calculation of the Constant K			
RID	RC	D	RC* 1/D
1	0.635	4.264	0.149
2	0.635	8	0.079
3	0.635	4.264	0.149
4	1.299	10	0.130

5	0.635	5	0.127
6	1.299	7	0.186
7	1.299	4	0.325
8	1.299	7	0.186
9	0.635	5	0.127
10	1.299	2	0.649
11	0.635	4	0.159
12	1.299	3.891	0.334
13	0.635	3	0.212
14	1.299	2	0.649
15	1.299	3	0.433
16	0.635	3	0.212
17	1.299	7	0.186
18	1.299	3	0.433
19	0.635	5	0.127
20	1.299	8	0.162
Sum			5.012
N			20
K (= N/Sum)			3.991

Note that the calculation includes data imputation, in which respondents with missing degree data (i.e., respondents 3 and 12) are assigned the estimated degree for their group. This is a procedure that was implicit in earlier RDS analyses (e.g., Heckathorn 2002, and Salganik and Heckathorn 2004), where degree estimates were made employing the available degree data and then assigned to each element in the recruitment table, irrespective of whether a recruit had missing degree data. Note also that data imputation extends to the seed (respondent 1).

A more conservative approach, which would entail loss of data, would have been to delete from the list any respondents for whom degree data are missing. However, the approach taken in this paper is to employ as much information as is validly available for calculating each of the two terms from which the population estimate is based, and then extend those estimators to the remainder of the data set. Thus, all available recruitment data are employed in the calculation of transition probabilities, and all available degree

data from peer-recruited respondents are employed to estimate degree.

The dual-component weight can now be calculated based on the recruitment and degree components. Respondent  $i$ 's dual-component weight,  $DW^i$ , is

$$(A8) \quad DW^i = RC^i \cdot DC^i = RC^i \cdot \frac{1}{D^i} \cdot K$$

The calculations of the degree component, and the dual-component sampling weight are illustrated in Table A.VI.

Table A.VI: Calculation of the Dual-Component Weight						
RID	RC	D	K	DC (=1/D*K)	DW (=RC*DC)	
1	0.635	4.264	3.991	0.936	0.594	
2	0.635	8	3.991	0.499	0.317	
3	0.635	4.264	3.991	0.936	0.594	
4	1.299	10	3.991	0.399	0.518	
5	0.635	5	3.991	0.798	0.507	
6	1.299	7	3.991	0.570	0.740	
7	1.299	4	3.991	0.998	1.296	
8	1.299	7	3.991	0.570	0.740	
9	0.635	5	3.991	0.798	0.507	
10	1.299	2	3.991	1.995	2.591	
11	0.635	4	3.991	0.998	0.633	
12	1.299	3.891	3.991	1.026	1.332	
13	0.635	3	3.991	1.330	0.845	
14	1.299	2	3.991	1.995	2.591	
15	1.299	3	3.991	1.330	1.727	
16	0.635	3	3.991	1.330	0.845	
17	1.299	7	3.991	0.570	0.740	
18	1.299	3	3.991	1.330	1.727	
19	0.635	5	3.991	0.798	0.507	
20	1.299	8	3.991	0.499	0.648	

When the vector of individualized weights is imported into a statistical program, it yields population estimates equal to those produced by the standard weights. A difference, however, is notable. When the mean degree by group is estimated, the categorical weight yields a result that ignores within-group variation in degree (i.e.,  $\widehat{D}_A =$

4.714,  $\widehat{D}_B = 5.3$ ), whereas when the individualized weight is used, the result coincides with the multiplicity-based degree estimate (i.e., 4.264 and 3.891 for groups A and B, respectively).

***Calculating the Recruitment Component for Degree***

The degree variable is partitioned into an appropriate aggregation level. For example, to simplify the analysis, it is divided into two categories, degree 2-5, and degree 6-10.

*Table A. VII: Recruitment Matrix by Degree*

---

Recruitment Count (Recruitment Proportion)	Recruits		Total
	Group Degree 2-5	Group Degree 6-10	
Person who Recruited			
Group Degree 2-5	7 (0.875)	1 (0.125)	8 1
Group Degree 6-10	2 (0.3333)	4 (0.6667)	6 1
Sample Proportional Composition (including seeds)	0.6111	0.3889	1
Equilibrium Sample Distribution	0.7273	0.2727	1

---

The recruitment component is then calculated using equation 19. The result is a vector of recruitment components.

*Table A.VIII: Calculation of the Recruitment Component of Degree (RCD)*

---

RID	Group	Seed?	$\hat{E}$	C	RCD (= $\hat{E}/C$ if respondent not a seed)
1	Degree 6-10	Yes	0.2727	0.3889	0
2	Degree 6-10	No	0.2727	0.3889	0.7013
3	0	No	0	0	0
4	Degree 6-10	No	0.2727	0.3889	0.7013
5	Degree 2-5	No	0.7273	0.6111	1.1901
6	Degree 6-10	No	0.2727	0.3889	0.7013
7	Degree 2-5	No	0.7273	0.6111	1.1901
8	Degree 6-10	No	0.2727	0.3889	0.7013
9	Degree 2-5	No	0.7273	0.6111	1.1901

---

10	Degree 2-5	No	0.7273	0.6111	1.1901
11	Degree 2-5	No	0.7273	0.6111	1.1901
12	0	No	0	0	0
13	Degree 2-5	No	0.7273	0.6111	1.1901
14	Degree 2-5	No	0.7273	0.6111	1.1901
15	Degree 2-5	No	0.7273	0.6111	1.1901
16	Degree 2-5	No	0.7273	0.6111	1.1901
17	Degree 6-10	No	0.2727	0.3889	0.7013
18	Degree 2-5	No	0.7273	0.6111	1.1901
19	Degree 2-5	No	0.7273	0.6111	1.1901
20	Degree 6-10	No	0.2727	0.3889	0.7013

Note that consistent with the practice of excluding seeds from degree calculations, a respondent is assigned a value of zero for the degree recruitment component not only if degree data are missing, but also if the respondent is a seed. Note also that RCD is nonneutral; that is, its value differs from one. Consequently, differential recruitment by degree is present in this data set.

Using group A as an example, the adjusted degree estimate is calculated as shown in Table A.IX.

Table A.IX. Calculating Adjusted Degree, Group A

RID	RCD	Degree	RCD/Degree
2	0.7013	8	0.0877
5	1.1901	5	0.2380
9	1.1901	5	0.2380
11	1.1901	4	0.2975
13	1.1901	3	0.3967
16	1.1901	3	0.3967
19	1.1901	5	0.2380
Sum	7.8418		1.8926

Adjusted degree = 4.1434 (= 7.8418/1.8926)

***Calculating the Adjusted Population Estimate***

The adjusted population estimate,  $\widehat{AP}$ , is then calculated by substituting the

adjusted degrees into the RDS estimator equation, as follows:

$$(A9) \quad \widehat{AP}_A = \frac{\widehat{S}_{BA} \widehat{AD}_B}{\widehat{S}_{BA} \widehat{AD}_B + \widehat{S}_{AB} \widehat{AD}_A} = \frac{0.2 \cdot 3.4523}{0.2 \cdot 3.4523 + 0.5 \cdot 4.1434} = 0.25$$

Based on the adjusted degree estimates, the estimated proportion of those in group A changes from 0.267 to 0.25. In this way, bias resulting from differential recruitment by degree can be controlled.

### *Calculating the Adjusted Dual-Component Weights*

The adjusted dual-component weights are calculated by multiplying the dual-component weights by the ratio of the adjusted and original RDS population estimator.

This is illustrated in Table A. X.

<i>Table A. X: Calculating the Adjusted Dual-Component Weight</i>					
RID	Variable	DW	AP	P	ADW (= DW(AP/P))
1	A	0.594	0.250	0.267	0.556
2	A	0.317	0.250	0.267	0.297
3	A	0.594	0.250	0.267	0.556
4	B	0.518	0.750	0.733	0.530
5	A	0.507	0.250	0.267	0.475
6	B	0.740	0.750	0.733	0.757
7	B	1.296	0.750	0.733	1.326
8	B	0.740	0.750	0.733	0.757
9	A	0.507	0.250	0.267	0.475
10	B	2.591	0.750	0.733	2.651
11	A	0.633	0.250	0.267	0.593
12	B	1.332	0.750	0.733	1.363
13	A	0.845	0.250	0.267	0.791
14	B	2.591	0.750	0.733	2.651
15	B	1.727	0.750	0.733	1.767
16	A	0.845	0.250	0.267	0.791
17	B	0.740	0.750	0.733	0.757
18	B	1.727	0.750	0.733	1.767
19	A	0.507	0.250	0.267	0.475
20	B	0.648	0.750	0.733	0.663

When these weights are imported into a statistics program such as STATA, and the

estimated frequency of the variable is analyzed using these weights, the result is the adjusted population estimate, 0.25 and 0.75 for groups A and B, respectively.

## Appendix B: Extending the Analysis beyond Dichotomous Variables: Data

### Smoothing

This section illustrates the means for deriving population estimates for variables with three or more categories using a data smoothing procedure. The following is Table IIA's recruitment matrix by age:

	Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101	RB
Age 20-33	13	10	3	2	2	30
Age 34-42	15	17	12	9	3	56
Age 43-49	7	9	8	14	12	50
Age 50-58	7	9	17	13	14	60
Age 59-101	8	4	10	15	17	54
RO	50	49	50	53	48	250

From Table IIA, the selection proportions are,

	Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101
Age 20-33	0.433	0.333	0.1	0.067	0.067
Age 34-42	0.268	0.304	0.214	0.161	0.054
Age 43-49	0.14	0.18	0.16	0.28	0.24
Age 50-58	0.117	0.15	0.283	0.217	0.233
Age 59-101	0.148	0.074	0.185	0.278	0.315

From Table IIA the equilibrium vector is,

Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101
0.233	0.219	0.186	0.192	0.17

In the demographically adjusted recruitment matrix each cell is the product of three terms, the selection proportion, and the equilibrium and total recruitments in the system.



*Table B.IV. Demographically Adjusted Recruitment Matrix*

	Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101	RB
Age 20-33	25.284	19.449	5.835	3.89	3.89	58.348
Age 34-42	14.671	16.627	11.737	8.803	2.934	54.772
Age 43-49	6.503	8.361	7.432	13.006	11.148	46.45
Age 50-58	5.587	7.184	13.569	10.376	11.175	47.891
Age 59-101	6.302	3.151	7.878	11.816	13.392	42.539
RO	58.347	54.772	46.451	47.891	42.539	250

For example, cell (2, 3), the number of recruitments by those age 34-42 of those age 43-49 is the  $0.21429 * 0.21909 * 250 = 11.73696$ .

Averaging the cross recruitment counts then produces the smoothed recruitment matrix,

*Table B.V. Data-Smoothed Recruitment Matrix*

	Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101	RB
Age 20-33	25.284	17.06	6.169	4.739	5.096	58.348
Age 34-42	17.06	16.627	10.049	7.993	3.043	54.772
Age 43-49	6.169	10.049	7.432	13.288	9.513	46.451
Age 50-58	4.739	7.993	13.288	10.376	11.495	47.891
Age 59-101	5.096	3.043	9.513	11.495	13.392	42.539
RO	58.348	54.772	46.451	47.891	42.539	250

For example, the counts in cell (2, 3) and cell (3, 2) are now equal to the mean of each, i.e.,  $(11.737 + 8.361)/2 = 10.049$ .

This smoothed matrix is then employed to recalculate all terms that are dependent upon recruitment counts, e.g., the smoothed selection proportion matrix is,

*Table B.VI. Data-Smoothed Transition Probabilities*

	Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101
Age 20-33	0.433	0.292	0.106	0.081	0.087
Age 34-42	0.311	0.304	0.183	0.146	0.056
Age 43-49	0.133	0.216	0.16	0.286	0.205
Age 50-58	0.099	0.167	0.277	0.217	0.24
Age 59-101	0.12	0.072	0.224	0.27	0.315

All other recruitment-count-dependent terms are then derived in the standard manner.

Drawing the dual-component degree estimates from Table IIA and the smoothed selection proportions from Table B.VI, the smoothed population estimate is derived by solving the following system of equations,

$$\begin{aligned}
 &1 = \widehat{P}_1 + \widehat{P}_2 + \widehat{P}_3 + \widehat{P}_4 + \widehat{P}_5 \\
 &\widehat{P}_1 \cdot 82.81 \cdot 0.292 = \widehat{P}_2 \cdot 109.777 \cdot 0.311 \\
 (B1) \quad &\widehat{P}_1 \cdot 82.81 \cdot 0.106 = \widehat{P}_3 \cdot 110.549 \cdot 0.133 \\
 &\widehat{P}_1 \cdot 82.81 \cdot 0.081 = \widehat{P}_4 \cdot 98.569 \cdot 0.099 \\
 &\widehat{P}_1 \cdot 82.81 \cdot 0.087 = \widehat{P}_5 \cdot 186.183 \cdot 0.12
 \end{aligned}$$

Solving this system of linear equations yields the adjusted population estimate reported in Table IA,

<i>Table B.VII. Adjusted Population Estimate</i>					
Age 20-33	Age 34-42	Age 43-49	Age 50-58	Age 59-101	
0.301	0.213	0.18	0.208	0.098	

## **Appendix C: Proof of equivalence of the RDS population estimators derived from categorical and individualized weights**

The difference between the categorical and individualized estimation procedures lies not in the computations upon which each estimate is based. Rather, they differ in the order in which these calculations are performed. An element of any multiplicity estimate is summing across degree reciprocals. In the categorical weights, this procedure is embedded within the calculation of the estimated mean for each group (i.e., see equation 21). In contrast, in estimates based on the individualized weights, this procedure occurs when the individualized weights are summed; because each respondent's individualized weight contains his or her own degree reciprocal (i.e., see equation 42). The equivalence of the two means of estimation can be demonstrated by unpacking each into its constituent terms and then simplifying the resulting expression.

To minimize the complexity of the proof, the simplest system from which each indicator can be calculated is used. Assume a dichotomous variable is to be analyzed with disjoint groups X and Y, with cross-recruitment proportions,  $S_{XY}$  and  $S_{YX}$ ; and the number of respondents in groups X and Y are  $n_X$  and  $n_Y$ , respectively. Assume further that only two respondents in each group have valid degree data, for this is the minimum amount of data for which a multiplicity adjustment is possible. The degrees are  $D^{X1}$  and  $D^{X2}$  for X, and  $D^{Y1}$  and  $D^{Y2}$  for Y.

### **Estimate Based on Categorical Weights**

The weighted estimate of a population proportion for a dichotomous variable is given by the expression,

$$(C1) \quad \widehat{P}_X = \frac{\sum_{i=1}^{n_X} W^i}{\sum_{i=1}^{n_X} W^i + \sum_{j=1}^{n_Y} W^j}$$

Given categorical weights, the values are the same for each member of group X,  $W_X$ , and for each member of Y,  $W_Y$ . Consequently, the above expression reduces to,

$$(C2) \quad \widehat{P}_X = \frac{n_X W_X}{n_X W_X + n_Y W_Y}$$

In expanded form, the weight for a member of group X is,

$$(C3) \quad W_X = \left( \frac{D^{Y1} D^{Y2} D^{X2} n_Y S_{YX} + D^{Y1} D^{Y2} D^{X1} n_Y S_{YX}}{D^{Y1} D^{Y2} D^{X2} n_Y S_{YX} + D^{Y1} D^{Y2} D^{X1} n_Y S_{YX} + D^{X1} D^{X2} D^{Y2} n_X S_{XY} + D^{X1} D^{X2} D^{Y1} n_X S_{XY}} \right) \left( \frac{n_X + n_Y}{n_X} \right)$$

The weight for members of group Y is defined in the same manner. When the weights for both groups are substituted into equation C2 above, and simplified, the result is,

$$(C4) \quad \widehat{P}_X = \frac{D^{Y1} D^{Y2} D^{X2} n_Y S_{YX} + D^{Y1} D^{Y2} D^{X1} n_Y S_{YX}}{D^{Y1} D^{Y2} D^{X2} n_Y S_{YX} + D^{Y1} D^{Y2} D^{X1} n_Y S_{YX} + D^{X1} D^{X2} D^{Y2} n_X S_{XY} + D^{X1} D^{X2} D^{Y1} n_X S_{XY}}$$

### Estimate Based on Individualized Weights

The estimate of a population proportion using individualized weights is given by the following expression, where each individualized weight as expanded into its constituents.

That is, substituting equation 42 into equation C1, and given that all Xs have the recruitment component  $RC_X$  and Y's the recruitment component  $RC_Y$ , yields the following,

$$(C5) \quad \widehat{P}_X = \frac{\sum_{i=1}^{n_X} \left( K \frac{1}{D^i} RC_X \right)}{\sum_{i=1}^{n_X} \left( K \frac{1}{D^i} RC_X \right) + \sum_{j=1}^{n_Y} \left( K \frac{1}{D^j} RC_Y \right)}$$

This estimate is a function of three types of terms, the recruitment component for each

group, the system constant, and respondents' degrees. Given that only two respondents in each group have valid degree data, this expression can be expanded as follows,

$$(C6) \quad \widehat{P}_X = \frac{\left( K \frac{1}{D^{X1}} RC_X \right) + \left( K \frac{1}{D^{X2}} RC_X \right)}{\left( K \frac{1}{D^{X1}} RC_X \right) + \left( K \frac{1}{D^{X2}} RC_X \right) + \left( K \frac{1}{D^{Y1}} RC_Y \right) + \left( K \frac{1}{D^{Y2}} RC_Y \right)}$$

The recruitment component can also be expanded and simplified, e.g., for X,  $RC_X$  expands to,

$$(C7) \quad RC_X = \frac{S_{YX} (n_X + n_Y)}{(S_{XY} + S_{YX}) n_X}$$

Similarly, the K term can be expanded,

$$(C8) \quad K = \frac{(S_{XY} + S_{YX}) n_X D^{X1} D^{X2}}{(D^{X1} + D^{X2}) S_{YX}}$$

Substituting the above expressions for K and  $RC_X$  and  $RC_Y$  into BC6, and then simplifying, produces the maximally expanded expression,

$$(C9) \quad \widehat{P}_X = \frac{D^{Y1} D^{Y2} D^{X2} n_Y S_{YX} + D^{Y1} D^{Y2} D^{X1} n_Y S_{YX}}{D^{Y1} D^{Y2} D^{X2} n_Y S_{YX} + D^{Y1} D^{Y2} D^{X1} n_Y S_{YX} + D^{X1} D^{X2} D^{Y2} n_X S_{XY} + D^{X1} D^{X2} D^{Y1} n_X S_{XY}}$$

Note that this expression is identical to equation C4 above. Consequently, the categorical and the individualized weights yield the same population estimate. Numerical analysis confirms that the conclusion derived from this simple case extends to the general case of systems with greater amounts of degree data and larger numbers of categories.