

Schonlau, Matthias; Liebau, Elisabeth

Working Paper

Respondent driven sampling

RatSWD Research Note, No. 45

Provided in Cooperation with:

German Data Forum (RatSWD)

Suggested Citation: Schonlau, Matthias; Liebau, Elisabeth (2010) : Respondent driven sampling, RatSWD Research Note, No. 45, Rat für Sozial- und Wirtschaftsdaten (RatSWD), Berlin

This Version is available at:

<https://hdl.handle.net/10419/43639>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

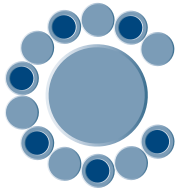
Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



German Data Forum
(RatSWD)

www.germandataforum.de

RatSWD

Research Notes

Research Note

No. 45

Respondent Driven Sampling

Matthias Schonlau, Elisabeth Liebau

September 2010

Research Notes of the German Data Forum (RatSWD)

The *RatSWD Research Notes* series publishes empirical research findings based on data accessible through the data infrastructure recommended by the RatSWD. The pre-print series was launched at the end of 2007 under the title *RatSWD Working Papers*.

The series publishes studies from all disciplines of the social and economic sciences. The *RatSWD Research Notes* provide insights into the diverse scientific applications of empirical data and statistics, and are thus aimed at interested empirical researchers as well as representatives of official data collection agencies and research infrastructure organizations.

The *RatSWD Research Notes* provide a central, internationally visible platform for publishing findings based on empirical data as well as conceptual ideas for survey design. The *RatSWD Research Notes* are non-exclusive, which means that there is nothing to prevent you from publishing your work in another venue as well: all papers can and should also appear in professionally, institutionally, and locally specialized journals. The *RatSWD Research Notes* are not available in bookstores but can be ordered online through the RatSWD.

In order to make the series more accessible to readers not fluent in German, the English section of the *RatSWD Research Notes* website presents only those papers published in English, while the German section lists the complete contents of all issues in the series in chronological order.

The views expressed in the *RatSWD Research Notes* are exclusively the opinions of their authors and not those of the RatSWD.

The *RatSWD Research Notes* are edited by:

Chair of the RatSWD (2007/ 2008 Heike Solga; 2009 Gert G. Wagner)

Managing Director of the RatSWD (Denis Huschka)

Respondent Driven Sampling

Matthias Schonlau*, Elisabeth Liebau**

**DIW Berlin, Germany and RAND Corporation, Pittsburgh, USA, (matt[at]rand.org)*

***DIW Berlin, Germany*

Abstract

Respondent driven sampling (RDS) is a network sampling technique typically employed for hard-to-reach populations (e.g. drug users, men who have sex with men, people with HIV). Similar to snowball sampling, initial seed respondents recruit additional respondents from their network of friends. The recruiting process repeats iteratively, thereby forming long referral chains. Unlike in snowball sampling, it is crucial to obtain estimates of respondents' personal network size (i.e., number of acquaintances in the target population) and information about who recruited whom. Markov chain theory makes it possible to derive population estimates and sampling weights. We introduce a new Stata program for RDS and illustrate its use.

JEL classification: C83, C88

Keywords: survey methodology, Stata software, chain referral sampling

1. Introduction

Some populations are difficult to sample. Consider the homeless: It is not possible to construct a sampling frame because there are no registries or other reasonably complete lists of the homeless. Random digit dialing does not work as most homeless are not known to carry around phones. Address based sampling procedures do not work well either because, well, the homeless do not have an address. Invented by Heckathorn in the mid-90ies, respondent driven sampling (RDS)(Heckathorn 1997, equation 6; Heckathorn 2002; Salganik and Heckathorn 2004) offers an alternative method that allows inference in populations for which traditional sampling methods are not feasible or not practical. RDS has proven particularly popular for behavioral surveillance of HIV and has been adopted by the Centers for Disease Control and Prevention (CDC).(Abdul-Quader et al. 2006)

Similar to snowball sampling, in RDS seed respondents recruit a fixed number of additional respondents from their network of friends. At each wave, recruits continue to recruit from among their friends. When the desired sample size is reached, the process is terminated. Unlike in snowball sampling, each respondent must be able to give an estimate of their network size (number of persons “you know” in the target populations; also called “degree”), and it is important to trace who recruited whom. Also unlike in snowball sampling, it is important that recruiting chains are sufficiently long to converge to a sampling equilibrium.

There are two additional features of RDS that the sampling theory does not require but which facilitate recruiting. First, there is a double incentive system. A respondent receives an incentive both for participating in the survey and for each successfully recruited respondent. Second, recruiting is driven by respondents rather than by interviewers. This feature also lends RDS its name. The idea is that respondents are more likely to participate when motivated by their friends, in particular when dealing

with a sensitive topic like AIDS or illegal drugs. In practice, respondents are only asked about their number of friends (“degree”), and some aggregate demographic information about their friends (“What percentage of your friends is female?”), but contact information of friends is never revealed to interviewers. Nevertheless, it is possible to track who recruited whom: respondents have to pass one coupon to each friend to be recruited. The coupon contains contact information about how to contact the interviewer or interviewing location and a unique coupon code. When the friend contacts the interviewer at a later time, the coupon code is collected as part of the meta-data. During analysis, the network structure can then be reconstructed – in fact, this is the purpose of *rds_network*, the first of our two programs.

The remainder of this article is organized as follows: Section 2 outlines some of the RDS theory including required assumptions. Section 3 contains information about the STATA implementation. Section 4 illustrates RDS by means of an example, the SATHCAP study, in detail. Section 5 concludes with a discussion.

2. Respondent Driven Sampling

Suppose we are interested in the population proportions of a categorical variable such as race/ethnicity or the prevalence of AIDS. We will call this variable an analysis variable and we will call each category (e.g. Hispanics) a group. Because we know who recruited whom, it is possible to compute a transition matrix of the analysis variable. RDS makes a Markov assumption: the value of the analysis variable of the recruited (e.g. Hispanic ethnicity) depends on the value of the analysis variable of the recruiter, but not on that of the recruiter’s recruiter.

For Markov chains transition matrix converges to a sample equilibrium and this equilibrium is independent of the seed (Heckathorn 2002, Theorem 1). Because of the independence to the seed, it

does not matter who the seed respondents are. (In practice, respondents thought to excel at recruiting, so-called social stars, are chosen as seed respondents.) The proportions in the sample equilibrium do not equal the population proportions, however, because respondents' inclusion probability is proportional to the number of their degree. That is, people who know more people in the target population are more likely to be recruited into the sample. Likewise, groups with larger average network size will be overrepresented in the equilibrium.

Estimating Average Group Degree

The network size of an individual respondent is called his or her degree. The average network size of a group is called average group degree. The average sample degree of a group is an overestimate of average group degree because respondents with a larger network are overrepresented in a sample. The multiplicity estimate of average degree (Heckathorn 2007, Section 2.1; Rothbart, Fine, and Sudman 1982) for group a corrects for this:

$$D_a = N_a / \sum_{i=1}^{N_a} (1/D_i)$$

where N_a is the sample size of group a and D_i is the degree of respondent i . (Seeds are excluded in the calculations of average group degree because seeds were not recruited by peers. (Heckathorn 2007, p.197; Salganik and Heckathorn 2004, p. 215)).

Estimating Population proportions

To derive population proportions, reciprocity or bi-directional recruiting relations are assumed. This means if respondent A recruited respondent B, then in principle the reverse could have occurred also. Denote k the total number of groups for which to compute population proportions, denote N_i , $i=1, \dots, k$ the sample sizes of group i . Further, denote S_{ij} the transition matrix between group i and j . Group i is the group of the recruiter and j the group of the recruit. The total number of ties originating from

members of group 1 is $N_1 D_1$, i.e. the number of respondents in group 1 times the average number of ties of group 1 respondents. The total number of ties between groups 1 and 2 can be computed as the total number of ties in group 1 times the proportion of ties that go from group 1 to group 2: $N_1 D_1 S_{12}$. Because of reciprocity, the total number of ties from group 2 to group 1, $N_2 D_2 S_{21}$, is equally large. Dividing by N turns the number of ties into population proportions, P_1 and P_2 , and the following equality is obtained (Heckathorn 2002, equation 8; Salganik and Heckathorn 2004, equation 6):

$$P_1 D_1 S_{12} = P_2 D_2 S_{21} \tag{1}$$

The constraint that proportions sum to 1 gives a second equation. If there are only two groups (e.g. is HIV positive or not), one can solve the two equations for the two unknown proportions. If there are more than two groups, equations analogous to (1) can be constructed for all pairs of groups. For m groups that yields $m*(m-1)/2$ equations (plus the constraint that proportions have to sum to 1) for only m parameters. The problem is over-determined. This dilemma can be solved, for example, by estimating the unknown parameters using least squares like in linear regression. Heckathorn's preferred solution, however, is a form of data smoothing (Heckathorn 2002, pp. 24-25). The underlying idea is follows: if groups recruit with equal effectiveness the number of people recruiting out of a group and into a group should be equal. The resulting demographically adjusted recruiting matrix R^* can be computed as follows (Heckathorn 2007, section 3.2):

$$R^* = \begin{pmatrix} S_{11}E_1N_r & S_{12}E_1N_r & \cdots & S_{1m}E_1N_r \\ S_{21}E_2N_r & S_{22}E_2N_r & \cdots & S_{2m}E_2N_r \\ \vdots & \vdots & \ddots & \vdots \\ S_{m1}E_mN_r & S_{m2}E_mN_r & \cdots & S_{mm}E_mN_r \end{pmatrix}$$

where N_r is the total number of recruits and $E_i, i=1, \dots, m$, is the proportion of group i in the equilibrium.

Because each row of the transition matrix is multiplied with a constant, $E_i * N_r$, the transition

probabilities are not affected. The smoothed demographically adjusted recruiting matrix R^{**} is a symmetric matrix where the smoothing consists of averaging:

$$R^{**} = \begin{pmatrix} R_{11}^* & \frac{R_{12}^* + R_{21}^*}{2} & \dots & \frac{R_{m1}^* + R_{1m}^*}{2} \\ \frac{R_{12}^* + R_{21}^*}{2} & R_{22}^* & \dots & \frac{R_{m2}^* + R_{2m}^*}{2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{R_{m1}^* + R_{1m}^*}{2} & \frac{R_{m2}^* + R_{2m}^*}{2} & \dots & R_{mm}^* \end{pmatrix}$$

Using the demographically adjusted recruiting matrix R^{**} , the transition matrix S^{**} can now be computed. Finally, proportion estimates can be obtained by solving the following system of m equations:

$$\begin{aligned} 1 &= P_1^{**} + P_2^{**} + \dots + P_m^{**} \\ P_1^{**} D_1 S_{12}^{**} &= P_2^{**} D_2 S_{21}^{**} \\ P_1^{**} D_1 S_{13}^{**} &= P_3^{**} D_3 S_{31}^{**} \\ &\vdots \\ P_1^{**} D_1 S_{1m}^{**} &= P_m^{**} D_m S_{m1}^{**} \end{aligned}$$

The smoothing renders additional equations redundant (Heckathorn 2007, p. 172). In case there are only 2 groups, the smoothing adjustment has no affect on the estimates of the proportions.

Sampling Weights

The population weights are computed by dividing the estimated population proportion for a given group equally among all sample members of that group:

$$W_a = P_a / C_a$$

where C_a refers to the sample proportion of group a (Heckathorn et al. 2002; Salganik and Heckathorn 2004). All members of group a receive the same population weight.

The population weight can be separated into a degree component, DC_a , and a recruitment component, RC_a , (Heckathorn 2007, equation 26):

$$W_a = (P_a / E_a) * (E_a / C_a) = DC_a * RC_a$$

The degree component represents a correction for differential average group degree. If the average group degrees are equal, then $P_a = E_a$ and the degree component $DC_a = 1$. The recruitment component represents differences in recruiting. When the sample proportion equals the equilibrium proportion, i.e. $C_a = E_a$, then the recruitment proportion $RC_a = 1$.

This partition leads to the introduction of individualized weights (Heckathorn 2007) or dual component weights DW_i .

$$DW_i = c * RC_i / D_i$$

where c is a normalizing constant chosen such that the average individualized weight equals 1.

Individualized weights are proportional to the inverse of a respondent's degree D_i . Because they are individualized, these weights are more appealing for individual level analyses.

Convergence

From theoretical work it is known that convergence to an equilibrium is reached fast (Heckathorn 2002, Theorem 2). Starting with an extreme distribution (100% of respondents in one group, 0% in all other groups), one can simulate how many recruitment waves are required for a given transition matrix to reach equilibrium. Convergence is achieved when two successive simulated recruitment waves do not

differ by more than a pre-specified convergence tolerance for any group. The Stata implementation of RDS requires that convergence is achieved from all m extreme distributions.

Homophily

Homophily measures to what extent respondents prefer to recruit from their own group rather than at random. The probability of selecting from the same group is the probability that selection is controlled by homophily plus the probability of random selection (Heckathorn 2002, p.20):

$$S_{aa} = H_a + (1 - H_a)P_a$$

for group a . Solving for H_a yields the equation for homophily. Homophily values range from -1 through +1. The value 0 corresponds to random recruitment; the value 1 corresponds to always recruiting from one's own group; the value -1 corresponds to never recruiting from one's own group. Moderate homophily is not problematic. If homophily is very large, however, the transition matrix may take a long time to converge which may be a sign that the groups are not networked.

The theory underlying RDS is based on a set of assumption which we explain in the following.

Assumption 1. Reciprocity. The reciprocity assumption implies that if respondent A recruited respondent B, then in principle B could have recruited A also. In practice, this assumption is tested by including a survey question about the relationship between the respondent and his or her recruiter. The assumption is violated if a lot of the recruited persons are strangers.

Assumption 2. Networked population. All respondents are interconnected. This assumption would be violated, for example, if the target population consisted of rivaling gangs who do not communicate with one another. The solution in this case would be to conduct separate RDS samples for

each of the non-communicating groups. If the number of waves required to reach an equilibrium for any variable is large, one may suspect a problem.

Assumption 3. Sampling with replacement. Sampling with replacement means that in principle a respondent could be contacted again and the respondent would participate a second time. In practice, a respondent would probably refuse to fill out the questionnaire a second time. In addition, duplicate respondents are usually actively screened out to prevent fraud related to obtaining multiple incentives. However, assuming that the sample is only a small fraction of the total population, this assumption can be ignored.

Assumption 4. Network size. Respondents can accurately report their personal network size. Biased estimates (e.g. consistent under- or overestimation of network size) are unproblematic as long as respondents uniformly under- or overestimate their network size (Wejnert 2009, Section "Degree Estimation"). There is ongoing concern that self-reported network sizes may be problematic (Wejnert and Heckathorn 2008, p.119), though there is also evidence that different ways of assessing network size lead to essentially the same result (Wejnert 2009).

Assumption 5. Random Recruitments. Respondents recruit from their network at random. To verify this assumption, one might ask about attributes (e.g. gender and race) of respondents' networks and compare expected characteristics to actual sample composition.

Nonresponse. Nonresponse matters in RDS also but is not talked about much. If respondents (rather than interviewers) recruit respondents, it is not possible to estimate nonresponse (unless the respondents are interviewed a second time). However, non-ignorable nonresponse would violate the random recruitment assumption and violations to this assumption can be tested.

3. Stata Implementation

RDS data look different than regular data because they embed the recruiting network structure. Table 1 gives an example of minimum data requirements: ID (coupon number), referral coupon numbers (here 6), network size, and an analysis variable (here race/ethnicity). The observations need not be ordered in any way. Missing referral coupons indicate that the respondents were not given a full set of referral coupons. In Table 1 no respondent was given more than 4 coupons. Whether a referral coupon was handed out but did not lead to a new respondent, or whether no coupon was handed out because sampling was terminated does not affect estimation.

The analysis is split into two Stata programs: *rds_network* and *rds*. The program *rds_network* determines the longest chain length (needed to assess convergence to the equilibrium), and it collects information about the recruiter of a respondent (variables *recruiter_id* and *recruiter_var*). The syntax is as follows:

```
rds_network varname , id(varname) coupon(str) ncoupon(int) degree(varname) ///  
[ ancestor(varname) depth(varname) recruiter_id(varname) recruiter_var(varname) ]
```

The options *id*, *coupon* and *ncoupon* specify the unique coupon code of respondents and their referral coupons, respectively. The program *rds_network* should always be called with the full RDS network for a given site. If a respondent is removed, the recruitment chain is broken into sub chains before and after the deleted respondent. (*Rds_network* intentionally does not support [if] and [in]). Optionally, the program generates two additional variables, *ancestor* and *depth*. *Ancestor* contains the id of the seed through which respondent was recruited. *Depth* contains the depth of the recruiting tree for a given recruit. Seeds have depth 0, their recruits have depth 1, and so forth.

rds is the main estimation program. The recruiter variables computed by *rds_network*, *recruiter_id* and *recruiter_var*, are now required as input variables. The syntax is as follows:

```
rds varname [if] [in] , id(varname) degree(varname) recruiter_id(varname) ///  
recruiter_var(varname) [ wgt(varname)]
```

Degree refers to the estimate of network size (number of friends in target population). Optionally, *wgt* generates a variable with individualized sampling weights. For clarity, some additional options (related to convergence to the equilibrium and the algorithm used to compute average network size) are not listed above.

Input validation and potential errors

The program *rds_network* verifies that respondent id and all referral coupons are unique. *rds_network* also verifies that there is no self-referral (a respondent's coupon points to him/herself). Further, *rds* will give an error if the estimated equilibrium proportion for a group is zero. Missing values for network size (*degree*) are allowed; missing values for the analysis variable specified in *<varname>* are not allowed. All network sizes (*degree*) must be positive.

Standard Errors and the Bootstrap

Standard errors and confidence intervals can be estimated via Taylor linearization (the *svy* routines in Stata) or by bootstrapping. The bootstrapping approach is preferred because of concern that the other approach does not adequately reflect variability in the sampling process. The bootstrap method is also the approach implemented in RDSAT (Volz et al. 2010). Even so, recent simulations suggest that confidence intervals are too typically too narrow (Goel and Salganik 2010). In Stata, *svy* routines can be applied as follows:

```
svyset [pweight=myweight]  
svy: proportion myvar
```

Standard errors of the proportions using a traditional nonparametric bootstrap of the ties between recruiter and recruitee are computed as follows:

```
bootstrap_b , reps(1000): rds varname, id() recruiter_id() [...]
```

The software RDSAT uses a slightly different bootstrapping procedure (Heckathorn 2002, pp. 27-29; Salganik 2006). Roughly, RDSAT simulates a new sample using the estimated transition matrix. The first simulated recruit is chosen arbitrarily. Each following simulated recruit is selected at random based on the probabilities specified in the transition matrix. In RDSAT, the bootstrapping procedure is applied to the “least squares” algorithm, not to the “smoothing” algorithm (RDSIncorporated 2006, p.30).

4. Example: SATHCAP study

The Sexual Acquisition and Transmission of HIV Cooperative Agreement Program (SATHCAP) applied RDS to sample men who had sex with men (MSM) and drug users (DU) in three US cities and in St. Petersburg, Russia (Iguchi et al. 2009). In addition, sex partners of this target population were sampled but were not part of the official RDS sample. The SATHCAP study used an innovative dual recruitment method with multi-colored coupons with different coupon colors for different segments of the target population to ensure both MSMs and DUs were sampled. The data used here to illustrate RDS corresponds to phase II at the Los Angeles site. We first analyze the network:

```
rds_network ethnic, id(id) coupon(numcpn) ncoupon(6) degree(netsize) recruiter_id(p_id) ///  
recruiter_var(p_key) depth(depth) ancestor(ancestor)
```

rds output (not shown) notes that there are 117 seed respondents.¹ This is an unusually large number.

The maximum chain or referral length is 18 (not counting the seed). The output also lists the length of

¹ The number of seeds reported in (Iguchi et al. 2009) is somewhat lower. During field work referral id’s of some respondents were lost. Rather than reporting the number of intended seeds, the program reports the number of actual seeds, namely respondents without a recruiter.

the maximal referral chains for each individual seed (Table 2 gives an excerpt). Most seeds in the Table 2 do not recruit anyone. Figure 1 shows the sample size by referral depth (using the variable *depth* specified above). Seeds have *depth=0*. The sample size decreases as the referral depth increases. Based on calculations with the variable specified in option *ancestor*, it turns out that 13 of the 117 seeds produce 71% of the sample. It is common that only a small percentage of seeds are highly productive.

Having computed the recruiter information, we can now proceed with assessing convergence and estimation:

```
rds ethnic, id(id) degree(netsize) recruiter_id(p_id) recruiter_var(p_key) wgt(wgt) wgt_pop(wgt2)
```

Originally, the variable *netsize* was calculated from 3 different questions corresponding to the number of MSMs, DUs and their overlap. Inconsistent answers could result in negative values and zeroes. We set those values to missing².

Convergence. The *rds* output (not shown) states that the required minimum referral length until convergence is 5. From the *rds_network* output we know that the longest chain in our data has length 18. Therefore, convergence for the variable “ethnic” is achieved. The required referral length needed to achieve convergence is simulated based on the transition matrix. It is also interesting to see how the sampling proportions converge. Figure 2 shows the sampling proportion of racial/ethnic groups calculated for all data up to a maximal depth or chain length. Indeed, we find that the proportions converge as the maximal wave increases, although in practice the convergence may have taken a little longer.

² Setting zeroes to 1 is less attractive as it would give those individuals very high weight. RDSAT routinely treats zeroes as missing.

Estimation. The final transition matrix is shown in Table 3. *Rds* output (not shown) contains intermediate matrices for the calculation of this transition matrix (S, R^*, R^{**}, S^{**}) and the matrix of observed counts. If there are only two groups the estimates of the initial and the final transition matrices are identical. In the transition matrix we notice that black respondents recruit other black respondents 67.5% of the time. We will get back to this in the context of homophily below.

Table 4 displays estimation results. The sample size is the sum of the number of seeds and the number of recruits. There were seeds in all four racial/ethnic categories. There are a total of three different proportion estimates: sample proportion, proportion in the equilibrium, and population proportion. The equilibrium proportion refers to the theoretical sampling proportion if the transition matrix has reached its equilibrium. If network size (degree) is constant, population proportions equal the equilibrium proportions. In practice, the network size varies and recruits who have a larger network are more likely to be sampled. The population proportion is an average-network-size-adjusted equilibrium proportion.

There are two measures of average network size in Table 4: “average” and “multiplicity”. The naïve estimate, “average” does not take into account that respondents with a larger network are more likely to be recruited into the sample. Therefore, the sample average for a group (e.g. Hispanics) overestimates the population average. The “multiplicity” estimate corrects for this. If the network sizes were constant then the two estimates would give the identical result.

The population sampling weights are designed to reproduce the estimated population proportion. The commands

```
svyset [pweight=wgt2]
```

```
svy: proportion ethnic
```

(where the variable *wgt2* was specified as an option in *rds*) reproduce the population proportions exactly. The variable weight contains only 4 distinct values corresponding to the four racial /ethnic categories.

Table 5 shows a comparison of estimated standard errors using Taylor linearization, bootstrap using RDS (estimates based on the “smoothing” algorithm introduced in section 2, n=2500) and the bootstrap from RDSAT (estimates based on the “least squares” algorithm, n=2500). The standard errors based on Taylor linearization are much smaller than the two bootstrapped estimates. The two bootstrap standard errors are similar to one another.

Homophily. Homophily is a diagnostic statistic that estimates to what extent respondents tend to recruit within-group rather than at random. For example, Table 4 shows that black respondents recruit 44.8% of the time other black respondents and 55.2% of the time they recruit at random from any of the 4 groups. Only very large homophily values would raise a concern.

Reciprocity. The SATHCAP questionnaire contained a question about the relationship between the respondent and his/her recruiter. It turns out only 4.5% of the recruited respondents described their recruiter as a stranger. This percentage is small and does not raise concerns. There are no guidelines of what percentage is considered too large or what to do if this assumption were violated.

Networked population. The number of iterations required to achieve convergence did not raise a red flag for any variable we looked at. Likewise, we found no anomalies in the corresponding transition matrices.

Random Recruitment. (Iguchi et al. 2009) argued it may not always be obvious to respondents how their friends self-identify in terms of race/ethnicity. Therefore, they looked at other variables including

gender to verify the random recruitment assumption. 88.7% of recruits are male (excluding a small number of transsexuals and excluding sex partners). Recruits reported 71.4% of their network is male. The difference is significant ($X^2(1)=74.0, p<0.001$). Therefore, the random recruitment assumption is violated with respect to gender. (Iguchi et al. 2009) argued that is not clear whether the differences are due to measurement error in the self-reported characteristics of their network or whether they are due to nonrandom recruitment.

5. Discussion

The integration of RDS within the Stata programming environment easily accommodates additional programming needs that require special purpose programming in a standalone package. For example, the bootstrap routine can be used with *rds* as explained earlier. Unusually large outliers of network size can be “pulled in” by setting large values to a user defined maximum. Further, some researchers might want to only analyze data after reaching equilibrium. If the equilibrium is reached after 5 referral waves, this can be accomplished as follows:

```
rds_network varname, depth(mydepth) [...]
rds varname if mydepth>=5, [...]
```

Weights can be poststratified to known totals using the *poststratify* option in *svyset* or, equivalently, a new adjusted weight variable can be computed using *svygen poststratify*.

There is currently no consensus on how to conduct regression with RDS data. Sampling weights are calculated based on a single analysis variable such as race/ethnicity. In multi-variable analyses such as regression it is unclear what to do. Current best practice is to conduct a sensitivity analysis (Johnston et al. 2008a) using the weight constructed for the dependent variable.

RDS is an area of active research and the literature is expanding. In practice, there are numerous implementation challenges such as defining eligibility criteria (Johnston 2008; Johnston et al. 2008b). Relative to a simple random sample, the RDS sample size should be at least twice as large to account for design effects and possibly larger (Goel and Salganik 2010; Salganik 2006). Recently, a second RDS estimator labeled RDSII has been derived (Volz and Heckathorn 2008). RDS has also been conducted through a web survey (Wejnert and Heckathorn 2008). We expect many more exciting developments on RDS in the future.

Acknowledgement

We are grateful to the SATHCAP project for allowing us to use some of their data in the example. Data from this project will be made publically available at a later date. We are grateful to Juergen Schupp whose intriguing and challenging questions inspired this project.

References

- Abdul-Quader, AS, DD Heckathorn, C McKnight, H Bramson, C Nemeth, K Sabin, K Gallagher, and DC 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 (3):459-476.
- Goel, S, and MJ Salganik. 2010. Assessing respondent-driven sampling. *Proceedings of the National Academy of Sciences* 107 (15):6743-6747.
- Heckathorn, DD. 1997. Respondent-driven sampling: a new approach to the study of hidden populations. *Social Problems* 44 (2):174-199.
- Heckathorn, DD. 2002. Respondent-driven sampling II: deriving valid population estimates from chain-referral samples of hidden populations. *Social Problems* 49 (1):11-34.
- Heckathorn, DD. 2007. Extensions of respondent-driven sampling: analyzing continuous variables and controlling for differential recruitment. *Sociological Methodology* 37 (1):151-207.
- Heckathorn, DD, S Semaan, RS Broadhead, and JJ Hughes. 2002. Extensions of respondent-driven sampling: a new approach to the study of injection drug users aged 18–25. *AIDS and Behavior* 6 (1):55-67.
- Iguchi, MY, AJ Ober, SH Berry, T Fain, DD Heckathorn, PM Gorbach, R Heimer, A Kozlov, LJ Ouellet, S Shoptaw, and Zule WA. 2009. Simultaneous Recruitment of Drug Users and Men Who Have Sex with Men in the United States and Russia Using Respondent-Driven Sampling: Sampling Methods and Implications. *Journal of Urban Health* 86 (Suppl 1):5-31.
- Johnston, L, H O’Bra, M Chopra, C Mathews, L Townsend, K Sabin, M Tomlinson, and C Kendall. 2008a. The associations of voluntary counseling and testing acceptance and the perceived likelihood of being HIV-infected among men with multiple sex partners in a South African township. *AIDS and Behavior* 14 (4):922-931.
- Johnston, LG. 2008. Behavioral Surveillance: Introduction to Respondent Driven Sampling (Participant Manual). Atlanta, GA: Centers for Disease Control and Prevention.
- Johnston, LG, M Malekinejad, C Kendall, IM Iuppa, and GW Rutherford. 2008b. Implementation challenges to using respondent-driven sampling methodology for HIV biological and behavioral surveillance: field experiences in international settings. *AIDS and Behavior* 12 (Suppl 1):131-141.
- RDSIncorporated. 2006. RDSAT 5.6 User Manual. Ithaca, NY.
- Rothbart, GS, M Fine, and S Sudman. 1982. On finding and interviewing the needles in the haystack: The use of multiplicity sampling. *Public Opinion Quarterly* 46 (3):408-421.
- Salganik, MJ. 2006. Variance estimation, design effects, and sample size calculations for respondent-driven sampling. *Journal of Urban Health* 83 (Suppl. 1):98-112.
- Salganik, MJ, and DD Heckathorn. 2004. Sampling and estimation in hidden populations using respondent-driven sampling. *Sociological Methodology* 34 (1):193-239.
- Volz, E, and DD Heckathorn. 2008. Probability based estimation theory for Respondent Driven Sampling. *Journal of Official Statistics* 24 (1):79-97.
- Volz, E, C Wejnert, I Degani, and DD Heckathorn. 2010. Respondent-Driven Sampling Analysis Tool (RDSAT) Version 6.0 Beta. Ithaca, NY: Cornell University.
- Wejnert, C. 2009. An empirical test of respondent-driven sampling: Point estimates, variance, degree measures, and out-of-equilibrium data. *Sociological Methodology* 39 (1):73-116.
- Wejnert, C, and DD Heckathorn. 2008. Web-based network sampling: Efficiency and efficacy of respondent-driven sampling for online research. *Sociological Methods & Research* 37 (1):105-134.

id	netsize	numcpn1	numcpn2	numcpn3	numcpn4	numcpn5	numcpn6	ethnic
40282	1	40307	40306	30306	30305	.	.	other
40361	3	40374	40375	30375	30376	.	.	white
172	18	40274	40275	other
40360	289	40383	40458	white
40383	12	30453	30454	40446	40447	.	.	black
40274	7	40335	40278	other
40275	4	40282	40283	other
40283	2	40361	40360	30359	30360	.	.	white
40278	6	40308	40309	white

Table 1: Example Data for RDS. The seed id appears in bold.

Seed	MaxDepth
...	...
2309	0
2378	0
2389	0
2395	0
2421	0
2462	2
2480	18
2499	1
2503	0
2602	0
...	...

Table 2: Excerpt of output from *rds_network* identifying seeds and the length of each seed's recruiting chain. Most seeds shown fail to recruit anyone.

	hispanic	white	black	other
hispanic	0.421	0.243	0.252	0.084
white	0.246	0.508	0.200	0.046
black	0.111	0.127	0.675	0.087
other	0.224	0.293	0.362	0.121

Table 3: Estimated Final Transition Matrix

	hispanic	white	black	other
Categories	1	2	3	4
SampleSize	160	167	282	55
Recruits	118	141	244	44
Seeds	42	26	38	11
Sample_Proportion	0.241	0.252	0.425	0.083
Equilibrium	0.226	0.268	0.427	0.078
PopulationProportion	0.260	0.249	0.411	0.079
AverageDegree	15.939	19.978	17.731	13.488
MultiplicityDegree	4.432	5.491	5.309	5.021
Homophily	0.217	0.344	0.448	0.045
Weight	1.081	0.992	0.967	0.959
RecruitmentComponent	0.939	1.067	1.006	0.943
DegreeComponent	1.151	0.929	0.961	1.016

Table 4: Estimation Results

	Taylor linearized std err	Bootstrap std err	RDSAT Bootstrap std err
hispanic	0.018	0.033	0.036
white	0.017	0.033	0.033
black	0.019	0.041	0.042
other	0.010	0.017	0.019

Table 5: Three estimates of the standard error of the population proportions of ethnicity: (1) Standard error based on Taylor approximation (using svyset), (2) bootstrap standard error (n=2500) using rds in stata, (3) bootstrap standard error (n=2500) using the RDSAT software.

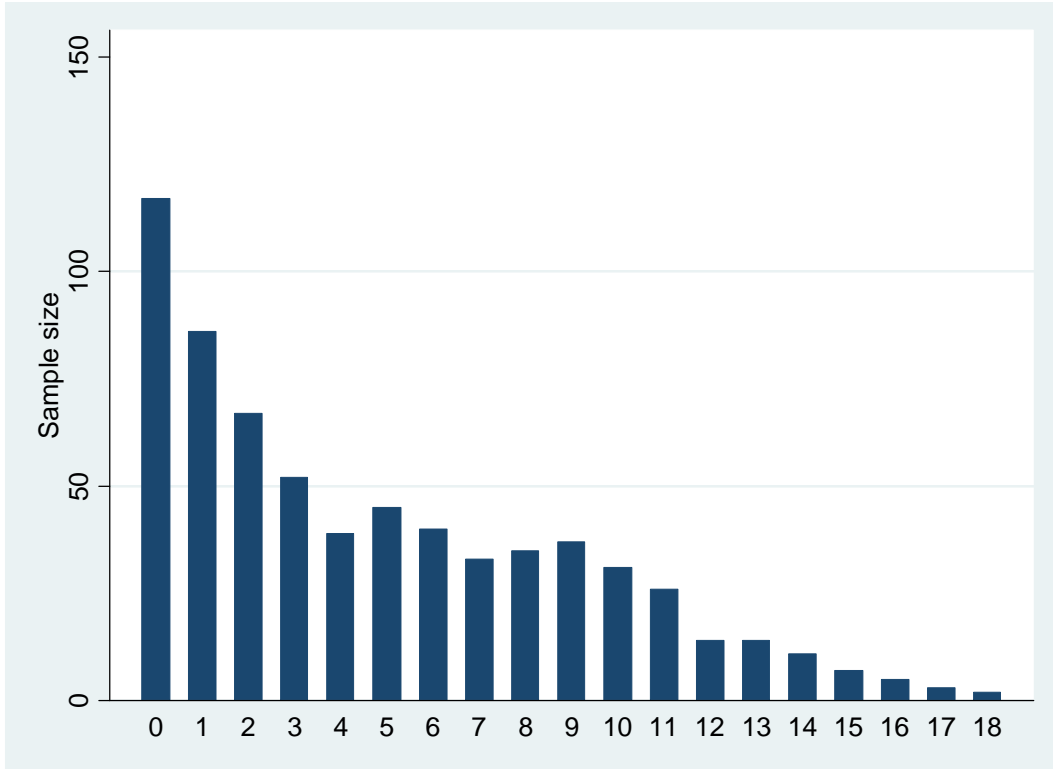


Figure 1: Sample size (excluding sex partners) by depth of the referral chain. Depth “0” corresponds to seed respondents.

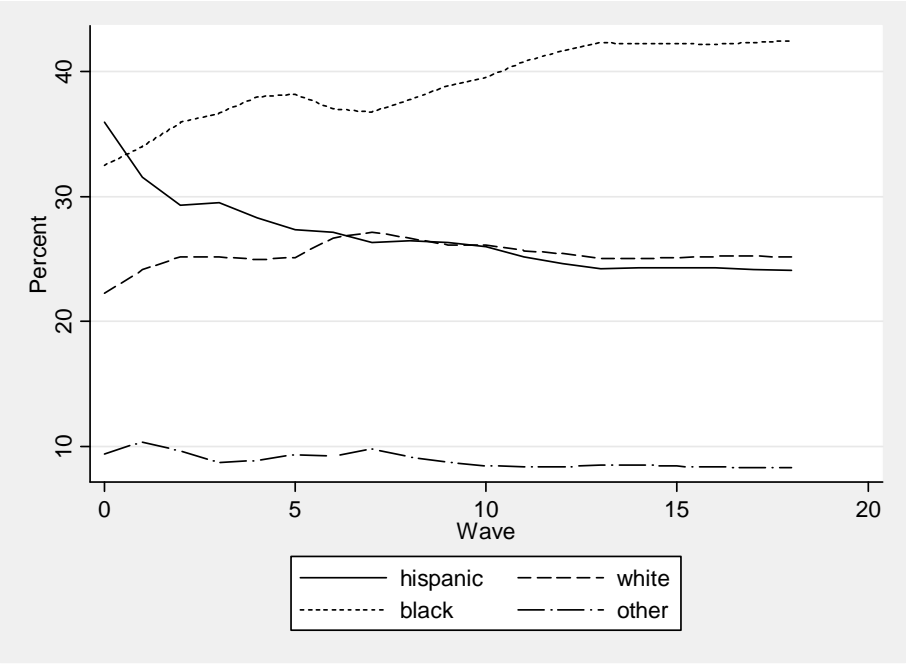


Figure 2: Percentage of four racial/ethnic groups for increasing length of the recruitment chain. Percentages are based on cumulative samples up to a given chain length.