

# Survival Analysis Estimation of an Eligibility Rate of Sampling Units Whose Eligibility Statuses Are Unknown: A Simulation Study of the Bias

Hiroaki Minato  
NORC at the University of Chicago

## Abstract

Our goal is to find the best estimator for the eligibility rate of the sampling units whose eligibility statuses are unknown. In this project, we focus on evaluating the survival analysis method proposed by Brick, Montaquila, and Scheuren (2000, 2002) in comparison to the method suggested by the Council of American Survey Research Organizations (CASRO) in 1982. The survival-analysis-based method is relatively new so its behavior is not yet completely understood. On the other hand, the CASRO method is typically thought of as the survey research industry's standard method. We compare the two estimation methods in terms of bias via simulation. We conclude that the bias in the survival-analysis-based method tends to be smaller than that in the CASRO method.

**Keywords:** Call history, Eligibility rate, Response rate, Survival analysis, Working residential number

## 1. Introduction

Our goal is to find the best estimator for the eligibility rate of the sampling units whose eligibility statuses are unknown. In this project, we focus on evaluating the survival analysis method proposed by Brick, Montaquila, and Scheuren (2000, 2002) in comparison to the method suggested by the Council of American Survey Research Organizations (CASRO) in 1982. Since the survival-analysis-based method is relatively new, its behavior is not yet completely understood. On the other hand, the CASRO method is typically thought of as the survey research industry's standard adjustment method. Using simulation techniques, we will compare the two methods with respect to bias within a simple but reasonable context.

It is important how and how well we estimate the eligibility rate for sampling units whose eligibility statuses are unknown; i.e., the estimation affects the response rate calculation. For a given sample, let  $E$  = the number of sampling units observed to be eligible,  $I$  = the number of sampling units observed to be ineligible, and  $U$  = the number of sampling units whose eligibility statuses could not be determined ( $E + I + U = n$ , some fixed total sample size). And, let  $u$  = the estimated eligibility rate of the sampling units of unknown eligibility status. Then, the estimated response rate is written as:

$$\hat{R} = \frac{C}{E + Uu},$$

where  $C$  is the number of eligible sampling units that responded ( $C \leq E$ ).

The true eligibility rate of the unknown eligibility status,  $\gamma$ , depends on the true nature of  $U$ . Since  $E + I + U = n$ ,  $U$  is

determined by  $E + I$  for a given  $n$ .  $C$  also depends on  $E$ .  $E + I$  are, in turn, a result of a particular data collection method and effort applied to a population of interest (e.g., sampling of telephone numbers and some calling rules associated with screening for interview eligibility) where  $E$  is defined by survey-specific requirements. Meanwhile, statistical properties of  $u$  are defined by a particular estimation method.

Given data  $C$ ,  $E$ , and  $U$ , we define the error in the estimated response rate as:

$$D = R - \hat{R} = \frac{C}{E + U\gamma} - \frac{C}{E + Uu} = \frac{CU(u - \gamma)}{(E + U\gamma)(E + Uu)}.$$

As we can see,  $u$  affects this error in a rather complicated way. (Trivially, however,  $D = 0$  for any  $u$ , if  $C = 0$  or  $U = 0$ .)

Consider a residential telephone survey. A typical set of calling rules generates eligible, ineligible, and undetermined sampling units. For each phone number sampled, calls are often attempted until the number is determined to be a working residential number (WRN) or otherwise (non-WRN); i.e., until eligibility is determined. However, since a maximum number of callbacks is often pre-specified, eligibility invariably fails to be resolved for some phone numbers. (We say calling rules are "full" calling rules, if all undetermined phone numbers are called until the maximum number is reached.)

Facing this situation, "in spite of many callbacks, eligibility status cannot be determined," CASRO (1982) suggests: "For purposes of estimating the number of eligible sampling units, this unknown remainder should be distributed between eligibles and ineligibles in the same proportion as exists among the working numbers." In other words, the CASRO method assumes cases of unknown eligibility are eligible WRN's in the same proportion as cases for which WRN status can be determined and it simply uses:

$$\frac{\text{Number of known WRN}}{\text{Number of known WRN} + \text{Number of known non-WRN}}$$

as the estimate of the WRN rate for the telephone numbers whose WRN statuses are unknown.

When the call history (call-level data) is available and when there is a variation in the WRN rate over call attempts, there is a promising alternative to the CASRO method. Brick, Montaquila, and Scheuren (2000, 2002) proposed the application of survival analysis in estimating the WRN rate among the undetermined cases. Their approach presupposes some random censoring of the number of call attempts for a given phone number, as the survival analysis estimate would be equivalent to the CASRO estimate without censoring. Thus, a proper comparison would be between the survival analysis method with *justifiable* "random-censoring" calling rules and the CASRO method with

“full” calling rules.

A topic-contributed session on the application of survival analysis method in analyzing call history data was organized at the 2004 Joint Statistical Meetings in Toronto. Luo and Minato (2005) theoretically compared, for some special cases, the bias in the survival analysis estimate versus the bias in the CASRO estimate. As we continue on that path, we provide some simulation results that indicate the circumstances under which the former bias is smaller than the latter bias.

## 2. Options for Fair Comparisons

If there are no numbers with undetermined eligibility at the end of eligibility screening (using “full” calling rules), there is nothing to do or to estimate. Also, if there are relatively few numbers left with undetermined eligibility, no additional attempts may be necessary or meaningful with respect to increasing the number of resolved or eligible cases. Nor would we need to apply a complicated estimation method for such cases. In fact, under these circumstances we might reasonably use the CASRO method or assume that all such cases are ineligible (or even assume that they are eligible, to be most conservative in calculating the response rate). In other words, if  $U$  is small, the impact of  $u$  on  $D$  would be relatively small.

Even when the number of resolved or eligible cases is not sufficiently large, if the resolution rate has decreased almost to zero by the end of screening, then there is little hope of increasing the number of such cases by accumulating further call attempts. In this situation, we might just stop screening and resort to using the CASRO method for estimating the eligibility rate among the undetermined cases. Note, in general, that when we have a large number of undetermined cases that cannot easily be resolved, we cannot do much or do well in terms of estimation as we usually have little or no information about their eligibility rates.

The alternative method of survival analysis becomes a reasonable option when (1) the number of resolved (and eligible) cases is not sufficiently large after using “full” calling rules and (2) further call attempts could bring in the sufficient number of resolved (and eligible) cases necessary to make up for the shortage. To compare the survival analysis method and CASRO method fairly, we may also need to fix the budget in terms of the total number of phone calls. As mentioned before, if the number of undetermined cases is large, neither method is expected to do well in taming bias.

First, suppose that a population consists of three types of phone numbers: for the first type of phone number ( $t_1$ ), eligibility can be resolved by making one call; the second type ( $t_2$ ) requires exactly two calls for resolution; and the third type ( $t_3$ ), three calls. Let  $\rho_1$ ,  $\rho_2$ , and  $\rho_3$  be the resolution rates for the first, second, and third call attempts, respectively.  $\rho_1$  and  $\rho_2$  must be greater than 0 but less than 1 so that  $t_1$ ,  $t_2$ , and  $t_3$  are all non-empty.  $\rho_3$  is set to 1 so that the union of  $t_1$ ,  $t_2$ , and  $t_3$  is the entire population. The eligibility rates (WRN rates) associated with  $t_1$ ,  $t_2$ , and  $t_3$  are written as  $\omega_1$ ,  $\omega_2$ , and  $\omega_3$ , respectively, where

$\omega_i$  can take on any value between 0 and 1, inclusively.

Assuming eligibility resolution for a population of phone numbers as defined above, the “full” calling rule can be reduced to a call-every-number-only-once rule. Now, subsequent to the first round of calls, consider the following two options: (1) call all the unresolved numbers one more time, or (2) censor out some of the unresolved numbers (at the rate  $c_1$ ) and call the remaining numbers one more time.

If we let  $c_2$  be the censoring rate after the second round of calls, then Option 1 can be expressed by setting  $c_1 = 0$  and  $c_2 = 1$  while Option 2 can be expressed by setting  $0 < c_1 < 1$  and  $c_2 = 1$ . Note that  $t_3$  cannot be observed because we assume a third call attempt would never be made under Options 1 and 2. In fact, it is reasonable to assume that some part of the population could never be reached using a realistic calling rule. Here, the first round of call attempts represents the part of the population for which eligibility can be resolved using “full” calling rules, while the second round of call attempts represents the part of the population for which additional calling attempts might prove promising.

Note that under Option 2, because fewer calls are made, cost will be lower. However, consequently you will end up with fewer resolved cases. Thus, we present the four-call-attempt scenario and a third option (see Table 1) that allows for a “fair” comparison between the CASRO method and the survival analysis method. Specifically, by dividing the second call attempt group into two groups, Option 2 can be extended such that *two* extra call attempts are made, therefore fixing the total number of resolved cases to be the same for both Options 1 and 3.

Assuming a finite population ( $N_1 < \infty$ ), the total number of calls made for Option 1 would be  $N_1 + N_1(1 - \rho_1)$ , while the total numbers of calls for Option 3 would be  $N_1 + N_1(1 - \rho_1)(1 - c_{31}) + N_1(1 - \rho_1)(1 - c_{31})(1 - \rho_2)(1 - c_{32})$ . The total number of resolved cases for Option 1 would be  $N_1\rho_1 + N_1(1 - \rho_1)\rho_2$ , and the total number of resolved cases for Option 3 would be  $N_1\rho_1 + N_1(1 - \rho_1)(1 - c_{31})\rho_2 + N_1(1 - \rho_1)(1 - c_{31})(1 - \rho_2)(1 - c_{32})\rho_3$ . Thus, constraining the cost (i.e., the total number of calls) and the return (i.e., the total number of resolved cases), we might set:

$$N_1 + N_1(1 - \rho_1) = N_1 + N_1(1 - \rho_1)(1 - c_{31}) + N_1(1 - \rho_1)(1 - c_{31})(1 - \rho_2)(1 - c_{32})$$

and

$$N_1\rho_1 + N_1(1 - \rho_1)\rho_2 = N_1\rho_1 + N_1(1 - \rho_1)(1 - c_{31})\rho_2 + N_1(1 - \rho_1)(1 - c_{31})(1 - \rho_2)(1 - c_{32})\rho_3.$$

For Option 3,  $0 < \rho_1 < 1$ ,  $0 < \rho_2 < 1$ , and  $0 < \rho_3 < 1$ ;  $0 < c_{31} < 1$  and  $0 < c_{32} < 1$ . So, under these constraints, the solution to the system of equations is:

$$\rho_2 = \rho_3 \tag{1}$$

and

$$c_{32} = \frac{(1 - \rho_3)(1 - c_{31}) - c_{31}}{(1 - \rho_3)(1 - c_{31})} \tag{2}$$

Obviously,  $\rho_2 = \rho_3$  is a rather limiting and somewhat unrealistic condition. However, in order to fairly compare Options 1 and 3 constraining the cost and return, it is necessary to make such an assumption.

### 3. Estimating $\gamma$

In this section, we give the CASRO and survival analysis estimates of  $\gamma$  under Options 1 and 3, respectively. Since the undetermined cases are generated in different ways by the two calling options (even though each assumes the same number of calls), we have two parameters to consider. Let  $\gamma_1$  be the true eligibility among the undetermined cases under Option 1, and let  $\gamma_3$  be the true eligibility rate under Option 3. Note that under Option 3 with up to four call attempts,  $N_1 > 0$ ,  $0 < \rho_1 < 1$ ,  $0 < \rho_2 < 1$ ,  $0 < \rho_3 < 1$ ,  $\rho_4 = 1$ ,  $0 \leq \omega_1 \leq 1$ ,  $0 \leq \omega_2 \leq 1$ ,  $0 \leq \omega_3 \leq 1$ , and  $0 \leq \omega_4 \leq 1$ , we have

$$\gamma_1 = \frac{N_3 \rho_3 \omega_3 + N_4 \rho_4 \omega_4}{N_3 \rho_3 + N_4 \rho_4}, \text{ where } N_3 = N_1(1 - \rho_1)(1 - \rho_2) \text{ and } N_4 = N_1(1 - \rho_1)(1 - \rho_2)(1 - \rho_3),$$

or

$$\gamma_1 = \frac{(1 - \rho_1)(1 - \rho_2)\rho_3\omega_3 + (1 - \rho_1)(1 - \rho_2)(1 - \rho_3)\omega_4}{(1 - \rho_1)(1 - \rho_2)\rho_3 + (1 - \rho_1)(1 - \rho_2)(1 - \rho_3)} \text{ with } \rho_4 = 1.$$

Further, with the censoring rates  $0 < c_1 < 1$ ,  $0 < c_2 < 1$ , and  $c_3 = 1$  (dropping the option index),

$$\gamma_3 = \frac{N'_2 \rho_2 \omega_2 + (N'_3 + N''_3) \rho_3 \omega_3 + (N'_4 + N''_4 + N'''_4) \rho_4 \omega_4}{N'_2 \rho_2 + (N'_3 + N''_3) \rho_3 + (N'_4 + N''_4 + N'''_4) \rho_4},$$

$$\text{where } N'_2 = N_1(1 - \rho_1)c_1, \quad N'_3 = N_1(1 - \rho_1)c_1(1 - \rho_2), \quad N''_3 = N_1(1 - \rho_1)(1 - c_1)(1 - \rho_2)c_2, \quad N'_4 = N_1(1 - \rho_1)c_1(1 - \rho_2)(1 - \rho_3), \quad N''_4 = N_1(1 - \rho_1)(1 - c_1)(1 - \rho_2)c_2(1 - \rho_3), \quad N'''_4 = N_1(1 - \rho_1)(1 - c_1)(1 - \rho_2)(1 - c_2)(1 - \rho_3)c_3, \text{ and}$$

or

$$\gamma_3 = [(1 - \rho_1)c_1\rho_2\omega_2 + \{(1 - \rho_1)c_1(1 - \rho_2) + (1 - \rho_1)(1 - c_1)(1 - \rho_2)c_2\}\rho_3\omega_3 + \{(1 - \rho_1)c_1(1 - \rho_2)(1 - \rho_3) + (1 - \rho_1)(1 - c_1)(1 - \rho_2)c_2(1 - \rho_3) + (1 - \rho_1)(1 - c_1)(1 - \rho_2)(1 - c_2)(1 - \rho_3)\}\omega_4] / [(1 - \rho_1)c_1\rho_2 + \{(1 - \rho_1)c_1(1 - \rho_2) + (1 - \rho_1)(1 - c_1)(1 - \rho_2)c_2\}\rho_3 + (1 - \rho_1)c_1(1 - \rho_2)(1 - \rho_3) + (1 - \rho_1)(1 - c_1)(1 - \rho_2)c_2(1 - \rho_3) + (1 - \rho_1)(1 - c_1)(1 - \rho_2)(1 - c_2)(1 - \rho_3)]$$

with  $\rho_4 = 1$  and  $c_3 = 1$ .

Letting  $n$ 's,  $r$ 's,  $w$ 's as the sample notation of  $N$ 's,  $\rho$ 's, and  $\omega$ 's, respectively, the CASRO estimate of  $\gamma_1$  is:

$$\hat{\gamma}_1 = \frac{n_1 r_1 w_1 + n_2 r_2 w_2}{n_1 r_1 + n_2 r_2}, \text{ where } n_2 = n_1(1 - r_1),$$

or

$$\hat{\gamma}_1 = \frac{r_1 w_1 + (1 - r_1) r_2 w_2}{r_1 + (1 - r_1) r_2}.$$

The survival analysis approach described by Brick, Montaquila, and Scheuren (2000, 2002) takes advantage of the relationship between the level of difficulty (number of call attempts) in reaching a household and the WRN rate (the eligibility rate). Provided that number of calls to a given phone number is censored randomly, phone numbers left with undetermined eligibility at the end of data collection can be considered right-censored observations. The survivor function for such data, which describes the probability of a number being resolved at each call attempt, can be partitioned into separate functions for WRN (eligible) and non-WRN (ineligible). Using similar notation to that in Brick et al. (2000, 2002), the mode-specific survivor functions are written as

$$\hat{S}_{WRN}(t) = \sum_{t' \geq t} \frac{d_{WRN}(t')}{n(t')} \hat{S}(t')$$

and

$$\hat{S}_{nonWRN}(t) = \sum_{t' \geq t} \frac{d_{nonWRN}(t')}{n(t')} \hat{S}(t'),$$

where  $d_{WRN}(t')$  and  $d_{nonWRN}(t')$  are the number of cases resolved to be WRN and nonWRN at the  $t'$ -th call attempt, respectively,  $n(t')$  is the number of cases available for the  $t'$ -th call attempt, and

$$\hat{S}(t) = \prod_{t' < t} \frac{n(t') - (d_{WRN}(t') + d_{nonWRN}(t'))}{n(t')}$$

is the Kaplan-Meier estimate of the marginal survivorship function. The overall WRN rate is then computed as

$$\hat{r}_\infty = \frac{\hat{S}_{WRN}(0)}{\hat{S}_{WRN}(0) + \hat{S}_{nonWRN}(0)}.$$

Finally, the WRN rate of the cases with undetermined eligibility is estimated as

$$\hat{r}_{UN} = \frac{\hat{r}_\infty \cdot n_{TOT} - n_{WRN}}{n_{UN}},$$

where  $n_{TOT}$  is the total number of cases,  $n_{WRN}$  is the number of cases resolved as WRN, and  $n_{UN}$  is the number with undetermined eligibility.

The survival analysis estimate of  $\gamma_3$  in our sample notation is:

$$\begin{aligned} \check{\gamma}_3 &= \left[ \frac{r_1 w_1 + (1 - r_1) r_2 w_2 + (1 - r_1)(1 - r_2) r_3 w_3}{1 - (1 - r_1)(1 - r_2)(1 - r_3)} n_1 \right. \\ &\quad \left. - \{r_1 w_1 + (1 - c_1)(1 - r_1) r_2 w_2 + (1 - c_1)(1 - c_2)(1 - r_1)(1 - r_2) r_3 w_3\} n_1 \right] \\ &\quad / \{ (1 - r_1) r_2 c_1 + (1 - r_1)(1 - r_2) - (1 - c_1)(1 - c_2)(1 - r_1)(1 - r_2) r_3 \} n_1 \\ &= \left[ \frac{r_1 w_1 + (1 - r_1) r_2 w_2 + (1 - r_1)(1 - r_2) r_3 w_3}{1 - (1 - r_1)(1 - r_2)(1 - r_3)} \right. \\ &\quad \left. - \{r_1 w_1 + (1 - c_1)(1 - r_1) r_2 w_2 + (1 - c_1)(1 - c_2)(1 - r_1)(1 - r_2) r_3 w_3\} \right] \\ &\quad / \{ (1 - r_1) r_2 c_1 + (1 - r_1)(1 - r_2) - (1 - c_1)(1 - c_2)(1 - r_1)(1 - r_2) r_3 \}. \end{aligned}$$

(Note that the observed rates of  $r_1$ ,  $r_2$ ,  $w_1$ , and  $w_2$  are *not* the same with those in  $\hat{\gamma}_1$  because the data are different.)

An inspection of the formulas for  $\gamma_1$ ,  $\hat{\gamma}_1$ ,  $\gamma_3$ , and  $\check{\gamma}_3$  suggests

that  $\tilde{\gamma}_3$  is a more reasonable estimate for  $\gamma_3$  than  $\tilde{\gamma}_1$  is for  $\gamma_1$ .  $\gamma_1$  contains  $\rho_1, \rho_2, \rho_3, \omega_3$ , and  $\omega_4$ , but  $\tilde{\gamma}_1$  consists of  $r_1, r_2, w_1$ , and  $w_2$ . That is,  $\tilde{\gamma}_1$  not only lacks  $r_3, w_3$ , and  $w_4$  but also contains  $w_1$  and  $w_2$  that are irrelevant in  $\gamma_1$ . On the other hand,  $\tilde{\gamma}_3$  is a function of  $r_1, r_2, r_3, w_1, w_2$ , and  $w_3$ , while  $\gamma_3$  is a function of  $\rho_1, \rho_2, \rho_3, \omega_2, \omega_3$ , and  $\omega_4$  (with  $c_1$  and  $c_2$  being constants in both functions). Thus,  $\tilde{\gamma}_3$  includes all relevant parameter estimates except for  $w_4$ , which no estimator can include as there is no data for  $\omega_4$  (unless  $\omega_4$  is assumed to be some function of  $\omega_1, \omega_2$ , or/and  $\omega_3$ ).  $w_1$  is found both in  $\tilde{\gamma}_1$  and  $\tilde{\gamma}_3$ , and this could therefore be a common source of bias.

#### 4. Simulations of Squared Asymptotic Biases

We conducted straightforward simulations for the above four-call-attempt scenario using calling and estimation Options 1 and 3. We specified the following sets of parameters values:

$$\begin{aligned} \rho_1 &= \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\} & , \\ \rho_2 &= \rho_3 = \{0.1, 0.2, 0.3, 0.4\} & , \\ \rho_4 &= 1 & , \\ \omega_1 &= \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\} & , \\ \omega_2 &= \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\} & , \\ \omega_3 &= \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\} & , \\ \omega_4 &= \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\} & , \\ c_1 &= \frac{(1-c_2)(1-\rho_2)}{(1-c_2)(1-\rho_2)+1} \quad (\text{from [2]},) & \text{ and} \\ c_2 &= \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}. \end{aligned}$$

Thus, there are 4,743,684 parameters combinations. Using the Maple (symbolic computation) software, we computed the squared asymptotic biases for each combination of parameters and counted (1) the number of combinations such that  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) > \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$ , (2) the number of combinations such that  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) = \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$ , and (3) the number of combinations such that  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) < \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$ . The distribution of the squared asymptotic bias is also computed for each estimator. The simulation results can be found in Table 2.

In this particular simulation, we can see that  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) > \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$  occupies about two thirds of the frequency distribution. If we use  $\tilde{\gamma}_3$  instead of  $\tilde{\gamma}_1$  when  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) > \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$ , then the reduction of the squared asymptotic bias would be about 47% ((549,978.2 - 291,456.9) / 549,978.2 100%). Meanwhile, if we use  $\tilde{\gamma}_1$  instead of  $\tilde{\gamma}_3$  when  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) < \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$ , then the reduction of the squared asymptotic bias would be about 34% ((141,917.4 - 93,121.9) / 141,917.4 100%).

If we always make the right choice of estimator, i.e., if we pick the estimator that gives us a smaller squared asymptotic bias, then the total squared asymptotic bias would be reduced by 40% ((647,964.2 - 389,442.9) / 647,964.2 100%), compared to

always using  $\tilde{\gamma}_1$ , and by 11% ((438,238.4 - 389,442.9) / 438,238.4 100%), compared to always using  $\tilde{\gamma}_3$ .

Next, we ask if there are any systematic patterns in the parameter combinations that characterize each of the outcomes

$$\begin{aligned} \text{Sq.Asy.Bias}(\tilde{\gamma}_1) &> \text{Sq.Asy.Bias}(\tilde{\gamma}_3) & , \\ \text{Sq.Asy.Bias}(\tilde{\gamma}_1) &= \text{Sq.Asy.Bias}(\tilde{\gamma}_3) & , \text{ and} \\ \text{Sq.Asy.Bias}(\tilde{\gamma}_1) &< \text{Sq.Asy.Bias}(\tilde{\gamma}_3). \end{aligned}$$

Studying the simulation results has led us to propose the following conjecture and theorem.

**Conjecture 1.**  $\text{Sq.Asy.Bias}(\tilde{\gamma}_1) > \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$  for all  $\omega_1, \omega_2, \omega_3$ , and  $\omega_4$  in  $[0, 1]$  such that

$$\begin{aligned} \{ \omega_1, \omega_2 \} &> \{ \omega_3, \omega_4 \} \text{ or } \{ \omega_1, \omega_2 \} < \{ \omega_3, \omega_4 \}, \\ \text{where } \{ \omega_1, \omega_2 \} &\text{ means } \omega_1 > \omega_2, \omega_1 = \omega_2, \text{ or } \omega_1 < \omega_2 \text{ and } \{ \omega_3, \omega_4 \} \\ &\text{ means } \omega_3 > \omega_4, \omega_3 = \omega_4, \text{ or } \omega_3 < \omega_4, \\ \text{whenever } N_1 > 0, & 0 < \rho_1 < 1, 0 < \rho_2 < 1, 0 < \rho_3 < 1, \rho_4 = 1, \\ 0 < c_1 < 1, & 0 < c_2 < 1, \text{ and } c_3 = 1 \text{ satisfy} \\ N_1(1-\rho_1)(1-c_1) &+ N_1(1-\rho_1)(1-c_1)(1-\rho_2)(1-c_2) = N_1(1-\rho_1) \\ \text{and} & \\ N_1(1-\rho_1)(1-c_1)\rho_2 &+ N_1(1-\rho_1)(1-c_1)(1-\rho_2)(1-c_2)\rho_3 \\ = N_1(1-\rho_1)\rho_2. \end{aligned}$$

This conjecture is a refinement of Conjecture 2 in Luo and Minato (2005, p. 3935). Notice the strict inequality in  $\{ \omega_1, \omega_2 \} > \{ \omega_3, \omega_4 \}$  or  $\{ \omega_1, \omega_2 \} < \{ \omega_3, \omega_4 \}$ . And, the strongly monotonic relations  $\omega_1 < \omega_2 < \omega_3 < \omega_4$  and  $\omega_1 > \omega_2 > \omega_3 > \omega_4$  meet the hypothesis.

Chart 1 shows simulation results using:

$\rho_2 = \rho_3 = 0.5, c_1 = 0.2, c_2 = 0.5$  ( $\Rightarrow$  the same numbers of total phone calls and the same numbers of resolved cases) and  $(\omega_1, \omega_2, \omega_3, \omega_4) = (0.2, 0.4, 0.6, 0.8), (0.8, 0.6, 0.4, 0.2), (0.1, 0.3, 0.7, 0.9)$ , and  $(0.9, 0.7, 0.3, 0.1)$ .  $\rho_1$  is varied on  $(0, 1)$ , and the plot is smoothed as if  $\rho_1$  is continuously varied. We can observe the following:

1. As Conjecture 1 suggests, the squared asymptotic bias under the CASRO option is larger than that under the survival analysis method option: green vs. red and blue vs. yellow.
2. The relationships between  $\rho_1$  and the squared asymptotic bias are identical for  $\omega_1 < \omega_2 < \omega_3 < \omega_4$  and  $\omega_1 > \omega_2 > \omega_3 > \omega_4$ . That is, there is some symmetry. This, however, may not be totally surprising, given the symmetry in Conjecture 1.
3. The squared asymptotic bias is larger when  $\omega$ 's are more spread out: green vs. blue and red vs. yellow.
4. As  $\rho_1$  increases, the squared asymptotic bias increases. Also note that as  $\rho_1$  increases, the number of undetermined cases decreases and thus that the impact of the bias on the response rate computation might be attenuated. (Recall the formula for  $D$ .)

When the eligibility rate is uniform in the population, we have the following theorem.

**Theorem 1.**  $\text{Sq.Asy.Bias}(\hat{\gamma}_1) = \text{Sq.Asy.Bias}(\tilde{\gamma}_3)$  for all  $\omega_1, \omega_2, \omega_3$ , and  $\omega_4$  in  $[0, 1]$  such that

$$\omega_1 = \omega_2 = \omega_3 = \omega_4,$$

whenever  $N_1 > 0$ ,  $0 < \rho_1 < 1$ ,  $0 < \rho_2 < 1$ ,  $0 < \rho_3 < 1$ ,  $\rho_4 = 1$ ,  $0 < c_1 < 1$ ,  $0 < c_2 < 1$ , and  $c_3 = 1$  satisfy

$$N_1(1 - \rho_1)(1 - c_1) + N_1(1 - \rho_1)(1 - c_1)(1 - \rho_2)(1 - c_2) = N_1(1 - \rho_1)$$

and

$$\begin{aligned} & N_1(1 - \rho_1)(1 - c_1)\rho_2 + N_1(1 - \rho_1)(1 - c_1)(1 - \rho_2)(1 - c_2)\rho_3 \\ & = N_1(1 - \rho_1)\rho_2. \end{aligned}$$

**Proof.** Given the conditions, we can directly show with algebra that  $\text{Sq.Asy.Bias}(\hat{\gamma}_1) - \text{Sq.Asy.Bias}(\tilde{\gamma}_3) = 0$ . »«

This result is rather intuitive, because if the eligibility does not vary, the survival analysis method is not expected to gain or lose any more information than the CASRO method.

**Tables**

Table 1: Four-call-attempt scenario

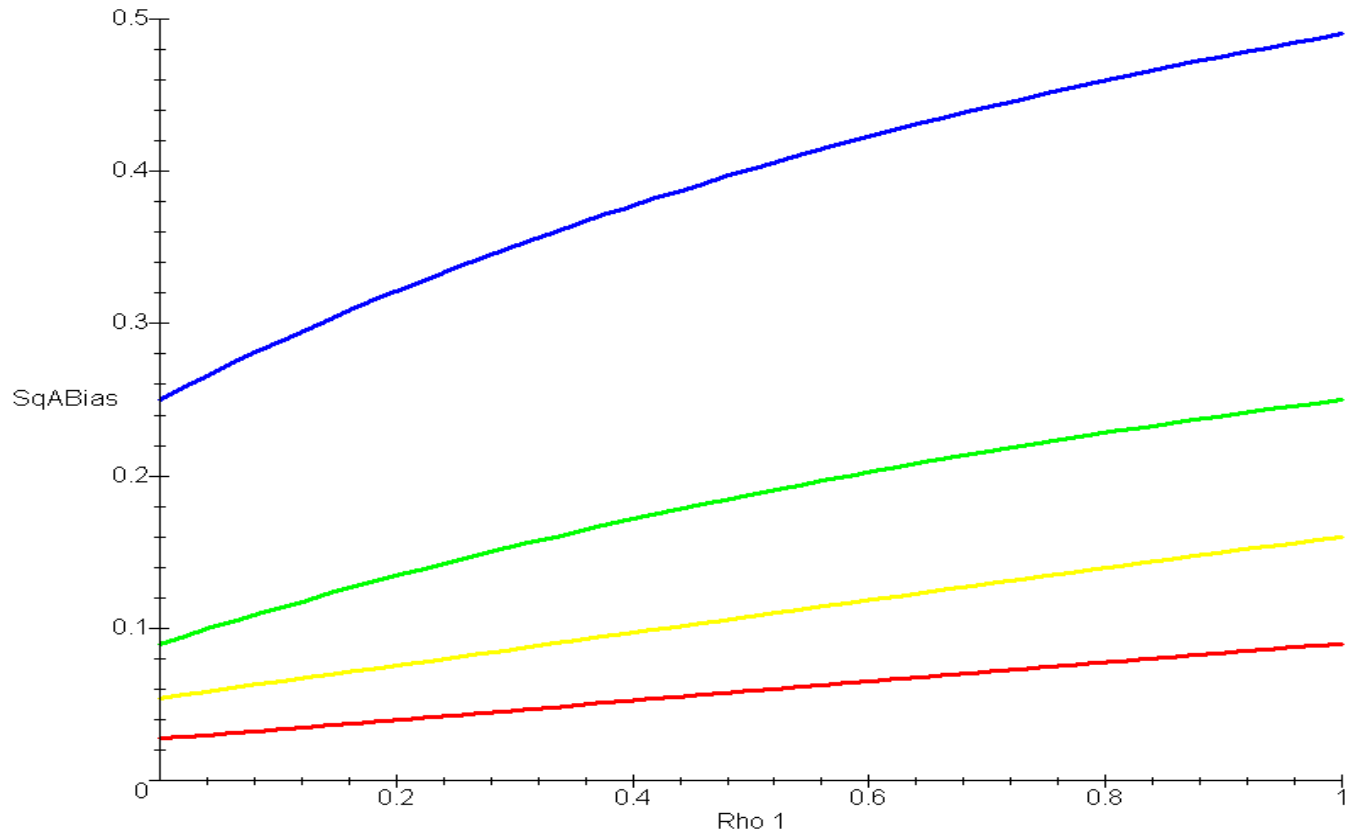
Call Attempt	Number of Calls	Resolution Rate	Resolved	Observed Eligibility Rate	Eligible	Non-eligible	Unresolved	Censoring Rate	Censored
1	$N_1$	$\rho_1$	$N_1\rho_1$	$\omega_1$	$N_1\rho_1\omega_1$	$N_1\rho_1(1-\omega_1)$	$N_1(1-\rho_1)$	$c_1$	$N_1(1-\rho_1)c_1$
2	$N_2 = N_1(1-\rho_1)(1-c_1)$	$\rho_2$	$N_2\rho_2$	$\omega_2$	$N_2\rho_2\omega_2$	$N_2\rho_2(1-\omega_2)$	$N_2(1-\rho_2)$	$c_2$	$N_2(1-\rho_2)c_2$
3	$N_3 = N_1(1-\rho_1)(1-c_1)(1-\rho_2)(1-c_2)$	$\rho_3$	$N_3\rho_3$	$\omega_3$	$N_3\rho_3\omega_3$	$N_3\rho_3(1-\omega_3)$	$N_3(1-\rho_3)$	1	$N_3(1-\rho_3)$
4		1		$\omega_4$					

Table 2: Simulation results

	Frequency (%)	Sum of Sq.Asy.Bias( $\hat{\gamma}_1$ ) (%)	Sum of Sq.Asy.Bias( $\tilde{\gamma}_3$ ) (%)	Sum of minimum of (Sq.Asy.Bias( $\hat{\gamma}_1$ ), Sq.Asy.Bias( $\tilde{\gamma}_3$ )) (%)
Sq.Asy.Bias( $\hat{\gamma}_1$ ) > Sq.Asy.Bias( $\tilde{\gamma}_3$ )	3,070,836 (64.7%)	549,978.2 (84.9%)	291,456.9 (66.5%)	291,456.9 (74.8%)
Sq.Asy.Bias( $\hat{\gamma}_1$ ) = Sq.Asy.Bias( $\tilde{\gamma}_3$ )	43,758 (0.9%)	4,864.1 (0.8%)	4,864.1 (1.1%)	4,864.1 (1.2%)
Sq.Asy.Bias( $\hat{\gamma}_1$ ) < Sq.Asy.Bias( $\tilde{\gamma}_3$ )	1,629,090 (34.3%)	93,121.9 (14.4%)	141,917.4 (32.4%)	93,121.9 (23.9%)
Total	4,743,684 (100%)	647,964.2 (100%)	438,238.4 (100%)	389,442.9 (100%)

Charts

Chart 1: Squared asymptotic bias



### Acknowledgements

The author thanks the Center for Excellence in Survey Research for the 2004 research grant that partly supported this research. The author also thanks Janella Chapline and Rachel Harter for reviewing the draft version of this paper.

### References

- The American Association for Public Opinion Research (2004), *Standard Definitions, Final Dispositions of Case Codes and Outcome Rates for Surveys: RDD Telephone Surveys, In-Person Household Surveys, Mail Surveys of Specifically Named Persons*, Lenexa, KS: Author.
- Andersen, P. K., Borgan, Ø., Gill, R. D., and Keiding, N. (1993), *Statistical Models Based on Counting Processes*, New York: Springer-Verlag.
- Brick, J. M., Montaquila, J., and Scheuren, F. (2000), "Estimating Residency Rates for Undetermined Telephone Numbers," in *American Statistical Association Proceedings of the Survey Research Methods Section*, pp. 1045-1050.
- (2002), "Estimating Residency Rates for Undetermined Telephone Numbers," *Public Opinion Quarterly*, 66, 18-39.
- Council of American Survey Research Organizations (1982), *On the Definition of Response Rates: A Special Report of the CASRO Task Force on Completion Rates*, Port Jefferson, NY: Author.
- Durand, C., Chevalier, S., and Vachon, S. (1998), "Using Survival Regression to Predict Occurrence of First Contact, Outcome at First Contact and Completion of Interview in Telephone Surveys," unpublished paper presented at the 14<sup>th</sup> World Congress of Sociology, Montreal, August, 1998.
- Lawless, J. F. (2003), *Statistical Models and Methods for Lifetime Data* (2<sup>nd</sup> ed.), Hoboken: John Wiley.
- Luo, L., and Minato, H. (2005). "Toward a Better Estimation of Working Residential Number (WRN) Rate Among the Undetermined: An Application of Survival Analysis," *2004 Proceedings of the American Statistical Association, Survey Research Methods Section* [CD-ROM], Alexandria, VA: American Statistical Association: 3933-3939.
- Smith, T. W. (2003), "A Review of Methods to Estimate the Status of Cases with Unknown Eligibility: Report Prepared for the AAPOR Standard Definitions Committee," unpublished paper presented at AAPOR, Phoenix, May, 2004.