

COMPARISON AND ANALYSIS OF STRATIFIED RANDOMIZED RESPONSE MODELS

Jong-Min Kim^a and Matthew E. Elam^b

^a Division of Science and Mathematics, University of Minnesota, Morris, MN, 56267, USA

^b Department of Industrial Engineering, The University of Alabama, Tuscaloosa, AL, 35487, USA

KEY WORDS: Random Sampling; Sensitive Characteristics; Estimation of Proportion

Introduction

Surveys are a means by which responses to questions concerning certain topics may be obtained from a sample of individuals selected in some manner from a population of interest. Results from surveys are affected by two main sources of error. The first is sampling error that results from taking a sample instead of enumerating the whole population. The second type of error is non-sampling error that cannot be attributed to sample-to-sample variability. Non-sampling error has two different errors which are random error and nonrandom error. Random error, which results from a reduction in the reliability of measurements, can be minimized over repeated measurements. However, nonrandom error, which is bias in the survey data, is difficult to cancel out over repeated measurements.

Deming (1960) and Cochran (1977) have discussed the sources of non-sampling error and its effects on sampling estimates. The main sources of non-sampling error in any survey are non-response bias and response bias. Non-response bias arises from subjects' refusal to respond and response bias arises from giving incorrect responses. When open or direct surveys are about sensitive matters (e.g., gambling habits, addiction to drugs and other intoxicants, alcoholism, proneness to tax evasion, induced abortions, drunken driving, history of past involvement in crimes, and homosexuality), non-response bias and response bias become serious problems because people oftentimes do not wish to give correct information.

In order to reduce non-response and response bias, a survey technique different from open or direct surveys was needed that made people comfortable and encouraged truthful answers. Warner (1965) developed such an alternative survey technique that is called randomized response (RR) technique. Warner's RR survey technique is designed to eliminate evasive answer bias and keep the respondents' confidentiality.

Since Warner presented the RR technique, many variants of the Warner model have been presented in the literature to further improve the technique, especially in regard to increasing the cooperation of the respondents and decreasing the variances of the

RR estimators (i.e., improving model efficiency). This paper briefly introduces three of these: stratified random sampling using optimal allocation versions of Warner's RR model, the unrelated question RR model, and the two-stage RR model. These models are compared with their respective counterparts that use simple random sampling. Furthermore, the proposed models are compared to each other in terms of relative efficiency. It is shown that stratified random sampling using optimal allocation further minimizes the variability in the estimate of the proportion of people with the sensitive trait under question, resulting in a survey with more precise results.

Simple Random Sampling RR Models

In initiating the work on the RR technique, Warner (1965) presented a two related question model for estimating the proportion of people who possess a sensitive trait in a given population. To apply the Warner model, a simple random sample of n people is drawn with replacement from the population. Before interviewing each person in the sample, each interviewer is furnished with an identical spinner which points to a statement 1 (I belong to the sensitive trait group) with probability P and to a statement 2 (I do not belong to the sensitive trait group) with probability $1 - P$. The interviewee spins the spinner unobserved by the interviewer. Without reporting the outcome of the spinner to the interviewer, the interviewee only answers "Yes" or "No" depending on the outcome of the randomization device. Warner (1965) equated the proportion of respondents who answer "Yes" to statement 1 or to statement 2 as follows:

$$X = P\pi_S + (1 - P)(1 - \pi_S) \quad (1)$$

where X is the proportion of "Yes" responses and π_S is the proportion of people with the sensitive trait. Warner (1965) derived several results from equation (1) and used these results to show his RR technique was an improvement over the regular methods available at the time in terms of reducing the variance in the estimate of π_S .

Greenberg et al. (1969) developed the theoretical framework for the unrelated question RR model suggested by Horvitz et al. (1967). Contrary to

Warner's (1965) model, the unrelated question RR model has one question that asks about a very sensitive trait and a second question that asks about an innocuous (or non-sensitive) trait. Greenberg et al. (1969) proposed two models, one for the case of an unknown π_I (the proportion of people with an innocuous trait) and the other for the case of a known π_I . For the case of an unknown π_I , simple random sampling with replacement is used to obtain two independent, nonoverlapping samples of sizes n_1 and n_2 from the population. Each interviewee in the i th sample is required to use Warner's (1965) randomization device with the outcomes having preassigned probabilities P_i and $1-P_i$ for $i=1, 2$. Without reporting the outcome of the spinner (statement A: I belong to the sensitive trait group; statement B: I belong to the innocuous trait group) to the interviewer, the interviewee answers "Yes" or "No" depending on the outcome from the randomization device. The proportion of respondents who answer "Yes" to statement A or to statement B is:

$$Y_i = P_i\pi_S + (1-P_i)\pi_I \text{ for } i=1, 2 \quad (2)$$

where Y_i is the proportion of "Yes" responses. Greenberg et al. (1969) derived several results from equation (2) and used these results to show their unrelated question RR technique was an improvement over Warner's (1965) RR technique in terms of reducing the variance in the estimate of π_S .

Mangat and Singh (1990) proposed a two-stage RR model that is a variation of Warner's (1965) model. In this model, each interviewee in the simple random sample with replacement of n respondents is provided with two randomization devices. The randomization device R_1 consists of two statements. The first statement is that "I belong to the sensitive trait group" (with probability M) and the second is "Go to randomization device R_2 " (with probability $1-M$). The randomization device R_2 also consists of two statements, which are "I belong to the sensitive group" and "I do not belong to the sensitive group" with known probabilities P and $1-P$, respectively. This is the same randomization device used by Warner (1965). Mangat and Singh (1990) derived the proportion of respondents who answer "Yes" to the sensitive question and to the negative of the sensitive question as:

$$\theta = M\pi_S + (1-M)[P\pi_S + (1-P)(1-\pi_S)] \quad (3)$$

where θ is the proportion of "Yes" responses. Mangat and Singh (1990) derived several results from equation (3) and used these results to show their two-stage RR technique was an improvement over Warner's (1965) RR technique in terms of reducing the variance in the estimate of π_S .

Stratified RR Models Using Optimal Allocation

In stratified random sampling, the population to be used to conduct the survey is partitioned into strata. A sample is then selected by simple random sampling with replacement from each stratum. To get the full benefit from stratification, it is assumed that the number of units in each stratum is known. In the stratified Warner's and unrelated question RR models, an individual respondent in the sample from stratum i is instructed to use the randomization device R_i which consists of a sensitive question (S) card with probability P_i and its negative question (\bar{S}) card with probability $1-P_i$. The respondent answers the question with a "Yes" or "No" without reporting which question card he or she has. A respondent belonging to the sample in different strata will perform different randomization devices, each having different preassigned probabilities. Under the assumption that these "Yes" and "No" reports are made truthfully and P_i is set by the researcher, the probability of a "Yes" answer in stratum i for the stratified Warner's RR model is:

$$Z_i = P_i\pi_{S_i} + (1-P_i)(1-\pi_{S_i}) \text{ for } i=1,2,\dots,k \quad (4)$$

where Z_i is the proportion of "Yes" answers in stratum i and π_{S_i} is the proportion of respondents with the sensitive trait in stratum i . Letting n_i denote the number of units in the sample from stratum i and n denote the total number of units in samples from all strata so that $n = \sum_{i=1}^k n_i$, Kim and

Warde (2003) give the optimal allocation of n to n_1, n_2, \dots, n_{k-1} and n_k to derive the minimum variance of $\hat{\pi}_S$ (an unbiased estimate of π_S) as follows:

$$\frac{n_i}{n} = \frac{w_i \left[\pi_{S_i}(1-\pi_{S_i}) + \frac{P_i(1-P_i)}{(2P_i-1)^2} \right]^{1/2}}{\sum_{i=1}^k w_i \left[\pi_{S_i}(1-\pi_{S_i}) + \frac{P_i(1-P_i)}{(2P_i-1)^2} \right]^{1/2}} \quad (5)$$

where $w_i = (N_i/N)$ for $i=1,2,\dots,k$ (N is the number of units in the whole population and N_i is the total number of units in stratum i). Under the assumptions that $n_i = n(N_i/N)$ and $P_i = P$ for all i , Kim and Warde (2003) give the minimal variance of $\hat{\pi}_S$ as follows:

$$\text{var}(\hat{\pi}_S) = \frac{1}{n} \left[\sum_{i=1}^k w_i \left\{ \pi_{S_i}(1-\pi_{S_i}) + \frac{P(1-P)}{(2P-1)^2} \right\}^{1/2} \right]^2 \quad (6)$$

For the stratified unrelated question RR model, the probability of a "Yes" answer in stratum i is:

$$Z_i = P_i\pi_{S_i} + (1-P_i)\pi_N \text{ for } i=1,2,\dots,k \quad (7)$$

where π_N , which is assumed known, is the proportion of respondents with the nonsensitive trait in stratum i . Kim and Elam (2003a) give the optimal allocation of n to n_1, n_2, \dots, n_{k-1} , and n_k to derive the minimum variance of $\hat{\pi}_S$ as follows:

$$\frac{n_i}{n} = \frac{\frac{w_i}{P_i} \sqrt{Z_i(1-Z_i)}}{\sum_{i=1}^k \frac{w_i}{P_i} \sqrt{Z_i(1-Z_i)}} \quad (8)$$

Kim and Elam (2003a) also give the minimal variance of $\hat{\pi}_S$ as follows:

$$\text{Var}(\hat{\pi}_S) = \frac{1}{n} \left[\sum_{i=1}^k \frac{w_i}{P_i} \sqrt{Z_i(1-Z_i)} \right]^2 \quad (9)$$

When π_N is unknown in the stratified unrelated question RR model, two independent non-overlapping simple random samples are drawn from each stratum. Two sets of the randomization device in each stratum are employed. The first set is used for respondents in the first sample, and the second set is used for respondents in the second sample. An individual respondent in the first sample from stratum i is instructed to use the randomization device R_{i1} which consists of a sensitive question (S) card with probability P_{i1} and a nonsensitive question (N) card with probability $1-P_{i1}$. An individual respondent in the second sample from stratum i is instructed to use the randomization device R_{i2} which consists of a

sensitive question (S) card with probability P_{i2} and a nonsensitive question (N) card with probability $1-P_{i2}$. The respondent answers the question with a "Yes" or "No" without reporting which question card he or she has in order to protect the respondent's privacy. So a respondent in different strata will perform different randomization devices, each having different preassigned probabilities. Let n_{i1} denote the number of units in the first sample from stratum i , n_{i2} denote the number of units in the second sample from stratum i , and n_i denote the total number of units in the two samples from each stratum. So $n = \sum_{i=1}^k n_i$ is the total number of units in the samples from every strata. Under the assumption that these "Yes" and "No" reports are made truthfully, the probability of a "Yes" answer in stratum i is:

$$Z_{i1} = P_{i1}\pi_{S_i} + (1-P_{i1})\pi_N$$

$$\text{and } Z_{i2} = P_{i2}\pi_{S_i} + (1-P_{i2})\pi_N \text{ for } i=1,2,\dots,k \quad (10)$$

where Z_{i1} is the proportion of "Yes" answers in the first sample from stratum i and Z_{i2} is the proportion of "Yes" answers in the second sample from stratum i . Kim and Elam (2003a) give the optimal allocation of n to n_1, n_2, \dots, n_{k-1} , and n_k to derive the minimum variance of $\hat{\pi}_S$ as follows:

$$\frac{n_i}{n} = \frac{\left[\frac{w_i \left((1-P_{i2})\sqrt{Z_{i1}(1-Z_{i1})} + (1-P_{i1})\sqrt{Z_{i2}(1-Z_{i2})} \right)}{(P_{i1}-P_{i2})} \right]}{\left[\sum_{i=1}^k \frac{w_i \left((1-P_{i2})\sqrt{Z_{i1}(1-Z_{i1})} + (1-P_{i1})\sqrt{Z_{i2}(1-Z_{i2})} \right)}{(P_{i1}-P_{i2})} \right]} \quad (11)$$

Kim and Elam (2003a) also give the minimal variance of the estimator $\hat{\pi}_S$ as follows:

$$\text{Var}(\hat{\pi}_S) = \frac{1}{n} \times \left[\sum_{i=1}^k \frac{w_i \left((1-P_{i2})\sqrt{Z_{i1}(1-Z_{i1})} + (1-P_{i1})\sqrt{Z_{i2}(1-Z_{i2})} \right)}{(P_{i1}-P_{i2})} \right]^2 \quad (12)$$

In the stratified two-stage RR model, the first stage of the survey interview requires an individual respondent in the sample from stratum i to use the randomization device R_{1i} which consists of a sensitive question (S) card with probability M_i and a “Go to the randomization device R_{2i} in the second stage” direction card with probability $1 - M_i$. The respondents in the second stage of stratum i are instructed to use the randomization device R_{2i} which consists of a sensitive question (S) card with probability P_i and its negative question (\bar{S}) card with probability $1 - P_i$. The respondent answers the question with a “Yes” or a “No” without reporting which question card he or she has in order to protect the respondent’s privacy. Under the assumption that M_i and P_i are set by the researcher, the probability of a “Yes” answer in stratum i is:

$$Y_i = M_i \pi_{S_i} + (1 - M_i)[P_i \pi_{S_i} + (1 - P_i)(1 - \pi_{S_i})]$$

for $i = 1, 2, \dots, k$ (13)

where Y_i is the proportion of “Yes” responses. Kim and Elam (2003b) give the optimal allocation of n to n_1, n_2, \dots, n_{k-1} and n_k to derive the minimum variance of $\hat{\pi}_S$ as follows:

$$\frac{n_i}{n} = \frac{w_i \left\{ \pi_{S_i} (1 - \pi_{S_i}) + \frac{(1 - M_i)(1 - P_i)[1 - (1 - M_i)(1 - P_i)]}{[2P_i - 1 + 2M_i(1 - P_i)]^2} \right\}^{1/2}}{\sum_{i=1}^k w_i \left\{ \pi_{S_i} (1 - \pi_{S_i}) + \frac{(1 - M_i)(1 - P_i)[1 - (1 - M_i)(1 - P_i)]}{[2P_i - 1 + 2M_i(1 - P_i)]^2} \right\}^{1/2}}$$

(14)

Kim and Elam (2003b) also give the minimal variance of the estimator $\hat{\pi}_S$ as follows:

$$\text{var}(\hat{\pi}_S) = \frac{1}{n} \left[\sum_{i=1}^k w_i \left\{ \pi_{S_i} (1 - \pi_{S_i}) + \frac{(1 - M_i)(1 - P_i)[1 - (1 - M_i)(1 - P_i)]}{[2P_i - 1 + 2M_i(1 - P_i)]^2} \right\}^{1/2} \right]^2$$

(15)

Comparison of RR Models

Kim and Warde (2003) use equation (6) to compare the efficiency of their stratified Warner’s RR model

to the Hong et al. (1994) stratified RR model using proportional sampling. Using Hong et al.’s (1994) $\text{var}(\hat{\pi}_H)$, Kim and Warde (2003) analytically show (see their equation (3.1.2)) that their model is more efficient when $P_i = P$ for all i . In the case that $P_i \neq P$ for all i , Kim and Warde (2003) perform an empirical analysis (see their Table 1) to show that their model is more efficient when there are two strata in the population, $P = P_1$, and $P_2 > P_1$.

Kim and Warde (2003) also use equation (6) to compare the efficiency of their model to the Mangat and Singh (1990) two-stage RR model. Using Mangat and Singh’s (1990) $\text{var}(\hat{\pi}_{ms})$, Kim and Warde (2003) show (see their Table 2) that their model is more efficient under the following condition when there are two strata in the population, $P = P_1 = P_2 \neq 0.5$, and $\pi_{S_1} \neq \pi_{S_2}$:

$$(\pi_{S_1} - \pi_{S_2})^2 + \left\{ \left[\pi_{S_1} (1 - \pi_{S_1}) + \frac{P(1 - P)}{(2P - 1)^2} \right]^{1/2} - \left[\pi_{S_2} (1 - \pi_{S_2}) + \frac{P(1 - P)}{(2P - 1)^2} \right]^{1/2} \right\}^2 > \left\{ \frac{M(1 - P)}{(2P - 1)(2P - 1 + 2M(1 - P))} \right\}^2 - \frac{M(1 - P)}{(1 - 2P)(2P - 1 + 2M(1 - P))^2} [w_1(1 - w_1)]^{-1}$$

(16)

By adding the additional condition $\pi_S > 1 - \{P/(2P - 1)\}^2$, Kim and Warde (2003) show (see their Theorem 3.2) that their model is more efficient than Mangat’s (1994) RR model when the following holds:

$$(\pi_{S_1} - \pi_{S_2})^2 + \left\{ \left[\pi_{S_1} (1 - \pi_{S_1}) + \frac{P(1 - P)}{(2P - 1)^2} \right]^{1/2} - \left[\pi_{S_2} (1 - \pi_{S_2}) + \frac{P(1 - P)}{(2P - 1)^2} \right]^{1/2} \right\}^2 > \frac{(1 - P)}{w_1(1 - w_1)P} \left[\left(\frac{P}{2P - 1} \right)^2 - \{1 - (w_1 \pi_{S_1} + w_2 \pi_{S_2})\} \right]$$

(17)

It should be noted that by setting $M = 0$ in the Mangat and Singh (1990) model (i.e., in equation

(3)), one gets Warner's (1965) RR model. For two strata in the population and $P = P_1 = P_2 \neq 0.5$, Kim and Warde's (2003) model is more efficient than Warner's (1965) model.

Kim and Elam (2003a) use equation (12) to compare the efficiency of their stratified unrelated question RR model when π_N is unknown to Kim and Warde's (2003) stratified Warner's RR model. Using equation (6) (with P_i replacing P), Kim and Elam (2003a) empirically show (see their Table 1) that their model is more efficient when $n=1000$, there are two strata in the population, $P_1 = P_{11} = P_{21}$, $P_2 = P_{12} = P_{22}$, and $P_1 + P_2 = 1$.

Kim and Elam (2003a) also use their equation (12) to compare the efficiency of their model when π_N is unknown to the Greenberg et al. (1969) unrelated question RR model. Using Greenberg et al.'s (1969) $\text{var}(\hat{\pi}_G)$, Kim and Elam (2003a) empirically show (see their Tables 2a and 2b) that their model is more efficient when $n=1000$, there are two strata in the population, $P_1 = P_{11} = P_{21}$, and $P_2 = P_{12} = P_{22}$.

Kim and Elam (2003b) use equation (15) to compare the efficiency of their stratified two-stage RR model to the Mangat and Singh (1990) two-stage RR model. Using Mangat and Singh's (1990) $\text{var}(\hat{\pi}_{ms})$, Kim and Elam (2003b) analytically show (see their Theorem 4.1) that their model is more efficient when there are two strata in the population, $P = P_1 = P_2 \neq 0.5$, and $M = M_1 = M_2$.

Kim and Elam (2003b) also use equation (15) to compare the efficiency of their model to Kim and Warde's (2003) stratified Warner's RR model. Using equation (6), Kim and Elam (2003b) analytically (see their Theorem 4.2) and empirically (see their Table 1) show that their model is more efficient under the condition $M > (1-2P)/(1-P)$ when there are two strata in the population, $P = P_1 = P_2 \neq 0.5$, $M = M_1 = M_2$, and $\pi_{S_1} \neq \pi_{S_2}$.

Discussion

This paper shows that stratified random sampling using optimal allocation further minimizes the variability in the estimate of the proportion of people with the sensitive trait under question for Warner's (1965) RR model, Greenberg et al.'s (1969) unrelated question RR model, and Mangat and Singh's (1990) two-stage RR model under the conditions presented. The result is a survey with more precise results. Additionally, as the number of strata increases, the variances calculated using equations (6), (12), and (15) decreases. This is shown in Section 3.3 of Kim and Warde (2003), Section 4.3 of Kim and Elam

(2003a), and Section 4.2 of Kim and Elam (2003b), respectively.

Two more advantages exist with stratified RR models using optimal allocation. The first is that they solve a limitation of RR which is the loss of individual characteristics of the respondents. Also, using optimal allocation helps to overcome the high cost incurred because of the difficulty in obtaining a proportional sample from a stratum (as in the Hong et al. (1994) model).

References

- Cochran, W.G., 1977. Sampling techniques, 3rd edn. New York: John Wiley and Sons.
- Deming, W.E., 1960. Sample design in business research. New York: John Wiley and Sons.
- Greenberg, B.G., Abul-Ela, A.A., Simmons, W.R., and Horvitz, D.G., 1969. The unrelated question randomized response: theoretical framework. *Journal of the American Statistical Association*, 64, 529-539.
- Hong, K., Yum, J., and Lee, H., 1994. A stratified randomized response technique. *Korean Journal of Applied Statistics*, 7, 141-147.
- Horvitz, D.G., Shah, B.V., and Simmons, W.R., 1967. The unrelated question randomized response model. *Proceedings of the Social Statistics Section of the American Statistical Association*, 65-72.
- Kim, J.-M. and Elam, M. E., 2003a. A stratified unrelated question randomized response model. *Journal of Statistical Planning and Inference*, in review.
- Kim, J.-M. and Elam, M. E., 2003b. A two-stage stratified Warner's randomized response model using optimal allocation. *Metrika*, in review.
- Kim, J.-M. and Warde, W.D., 2003. A stratified Warner's randomized response model. *Journal of Statistical Planning and Inference*, in press.
- Mangat, N.S. and Singh, R., 1990. An alternative randomized response procedure. *Biometrika*, 77, 439-442.
- Mangat, N.S., 1994. An improved randomized response strategy. *Journal of the Royal Statistical Society Series B*, 56, 93-95.
- Warner, S.L., 1965. Randomized response: a survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60, 63-69.