

Robert P. Clickner and Boris Iglewicz, Temple University

1. Introduction

Surveys, particularly those involving sensitive questions, have always been plagued by biases caused by non-response and untruthful responses. The first randomized response scheme designed to reduce these biases was developed by Warner [2]. Since then a number of generalizations and variations of Warner's technique have been developed. The recent survey paper by Horvitz, Greenberg and Abernathy [1] summarizes much of this work. The previous randomized response papers discuss the estimation of parameters and efficiencies of the randomized response techniques as compared to direct question surveys, but only in the context of analysis based on one sensitive question at a time. In actual surveys one is not only interested in the analysis of single characteristics, but also in joint estimates of a number of characteristics. Actually one is interested in obtaining as much useful information as possible and the popularity of 2x2 contingency tables reflects the need for cross-tabulation analysis of survey data.

In this paper we consider the extension of Warner's scheme to two sensitive questions. The maximum likelihood estimators of the joint and marginal probabilities are derived and their means, variances and covariances are obtained. Some properties of these estimators are explored and recommendations concerning choice of randomization parameters are made. Also considered are the efficiencies of these estimators relative to direct questioning; assuming both zero and positive probabilities of untruthful responses.

2. Statement of the Problem

The problem is to simultaneously estimate the proportion of the population who possess either or both of two sensitive characteristics using Warner's [2] randomized response scheme. If A and B are the sensitive characteristics let

$$\pi_A = P(\text{a person is an A}),$$

$$\pi_B = P(\text{a person is a B}),$$

$$\pi_{AB} = P(\text{a person is both an A and a B}).$$

We need only estimate π_A , π_B and π_{AB} . Once these three probabilities are estimated, all the other joint, marginal, and conditional probabilities concerning A and B can be easily estimated.

Warner's scheme, extended to the two question case, proceeds as follows. The respondent is given a random device (e.g., a spinner, deck of cards, box of colored marbles, etc.) with which to choose one of the two statements

1-a. I am an A,

1-b. I am not an A.

The device selects 1-a with probability p_1 and 1-b with probability $\bar{p}_1 = 1 - p_1$ (Given any probability θ , we will always use the notation

$\bar{\theta}$ for $1 - \theta$). Without revealing to the interviewer which statement has been chosen, the respondent answers "yes" or "no" according to the statement selected and to his actual status with respect to the characteristic A. This procedure is then repeated with the statements

2-a. I am a B,

2-b. I am not a B,

and with a second random device which selects 2-a with probability p_2 and 2-b with probability \bar{p}_2 . We need not have $p_1 = p_2$. Hence the information received from each respondent is one of the four pairs: yes-yes, yes-no, no-yes, no-no. We will code the responses as 1 = "yes" and 0 = "no."

Let $\lambda_{ij} = P(\text{response } i \text{ on the first question and response } j \text{ on the second question})$, $i = 0, 1$; $j = 0, 1$. Then writing \bar{A} for "not A" and \bar{B} for "not B," we obtain

$$\begin{aligned} \lambda_{11} &= p_1 p_2 P(AB) + p_1 \bar{p}_2 P(A\bar{B}) + \bar{p}_1 p_2 P(\bar{A}B) + \bar{p}_1 \bar{p}_2 P(\bar{A}\bar{B}) \\ &= p_1 p_2 \pi_{AB} + p_1 \bar{p}_2 (\pi_A - \pi_{AB}) + \bar{p}_1 p_2 (\pi_B - \pi_{AB}) \\ &\quad + \bar{p}_1 \bar{p}_2 (1 - \pi_A - \pi_B + \pi_{AB}) \\ &= \pi_A \bar{p}_2 (p_1 - \bar{p}_1) + \pi_B \bar{p}_1 (p_2 - \bar{p}_2) + \pi_{AB} (p_1 - \bar{p}_1)(p_2 - \bar{p}_2) \\ &\quad + \bar{p}_1 \bar{p}_2, \end{aligned} \quad (2.1)$$

$$\begin{aligned} \lambda_{10} &= \pi_A p_2 (p_1 - \bar{p}_1) - \pi_B \bar{p}_1 (p_2 - \bar{p}_2) - \pi_{AB} (p_1 - \bar{p}_1)(p_2 - \bar{p}_2) \\ &\quad + \bar{p}_1 p_2, \end{aligned} \quad (2.2)$$

$$\begin{aligned} \lambda_{01} &= -\pi_A \bar{p}_2 (p_1 - \bar{p}_1) + \pi_B p_1 (p_2 - \bar{p}_2) - \pi_{AB} (p_1 - \bar{p}_1)(p_2 - \bar{p}_2) \\ &\quad + p_1 \bar{p}_2 \end{aligned} \quad (2.3)$$

$$\begin{aligned} \lambda_{00} &= -\pi_A p_2 (p_1 - \bar{p}_1) - \pi_B p_1 (p_2 - \bar{p}_2) + \pi_{AB} (p_1 - \bar{p}_1)(p_2 - \bar{p}_2) \\ &\quad + p_1 p_2. \end{aligned} \quad (2.4)$$

3. Estimation of π_A , π_B , π_{AB} ; Test for Independence

Consider a sample of n respondents and let X_{ij} , $i = 0, 1$; $j = 0, 1$, be the number of respondents who respond i to question 1 and j to question 2. The joint distribution of X_{11} , X_{10} , X_{01} and X_{00} is multinomial $(n; \lambda_{11}, \lambda_{10}, \lambda_{01}, \lambda_{00})$. The maximum likelihood estimator (MLE) of λ_{ij} is $\hat{\lambda}_{ij} = X_{ij}/n$. Using equations (2.1)-(2.4) and some algebra one obtains the MLE's of π_A , π_B , π_{AB} to be

$$\hat{\pi}_A = (\hat{\lambda}_{11} + \hat{\lambda}_{10} - \bar{p}_1) / (p_1 - \bar{p}_1),$$

$$\hat{\pi}_B = (\hat{\lambda}_{11} + \hat{\lambda}_{01} - \bar{p}_2) / (p_2 - \bar{p}_2),$$

$$\begin{aligned} \hat{\pi}_{AB} &= [\hat{\lambda}_{11} (p_1 p_2 - \bar{p}_1 \bar{p}_2) - \hat{\lambda}_{10} \bar{p}_2 - \hat{\lambda}_{01} \bar{p}_1 + \bar{p}_1 \bar{p}_2] \\ &\quad / (p_1 - \bar{p}_1)(p_2 - \bar{p}_2), \end{aligned}$$

respectively. Note that $\hat{\pi}_A$ and $\hat{\pi}_B$ are exactly the same estimates as found in Warner [2].

The maximum likelihood estimates of π_A , π_B , π_{AB} are unbiased and have the following variances and covariances:

$$\begin{aligned} V(\hat{\pi}_A) &= [\pi_A \bar{\pi}_A + f(p_1)]/n, \\ V(\hat{\pi}_B) &= [\pi_B \bar{\pi}_B + f(p_2)]/n, \\ V(\hat{\pi}_{AB}) &= [\pi_{AB} \bar{\pi}_{AB} + \pi_A f(p_2) + \pi_B f(p_1) \\ &\quad + f(p_1)f(p_2)]/n, \\ \text{Cov}(\hat{\pi}_A, \hat{\pi}_B) &= (\pi_{AB} - \pi_A \pi_B)/n, \\ \text{Cov}(\hat{\pi}_A, \hat{\pi}_{AB}) &= [\pi_{AB} \bar{\pi}_A + \pi_B f(p_1)]/n, \\ \text{Cov}(\hat{\pi}_B, \hat{\pi}_{AB}) &= [\pi_{AB} \bar{\pi}_B + \pi_A f(p_2)]/n, \end{aligned}$$

where $f(p) = \bar{p}p/(p-\bar{p})^2$.

Let $\lambda_{i.} = \lambda_{i1} + \lambda_{i0}$, $i = 0, 1$ and $\lambda_{.j} = \lambda_{1j} + \lambda_{0j}$, $j = 0, 1$. Then we have $\lambda_{11} = \lambda_{1.} \lambda_{.1}$ if and only if $\pi_{AB} = \pi_A \pi_B$. That is, the responses to questions 1 and 2 are independent if and only if the characteristics A and B are independent. It follows that we can test the independence of A and B by applying the ordinary χ^2 test to the randomized responses. Specifically we can use the test statistic

$$\chi^2 = \sum_{i,j} (x_{ij} - n\hat{\lambda}_{i.}\hat{\lambda}_{.j})^2 / n\hat{\lambda}_{i.}\hat{\lambda}_{.j},$$

where $\hat{\lambda}_{i.} = \hat{\lambda}_{i1} + \hat{\lambda}_{i0}$ and $\hat{\lambda}_{.j} = \hat{\lambda}_{1j} + \hat{\lambda}_{0j}$.

4. Estimation of π_{AB} Only

Occasions may arise in which one is only interested in estimating π_{AB} . In such situations the question arises as to which is the better procedure--the two-question procedure described in Section 2 or Warner's original procedure applied to A and B, that is, asking the respondent to answer one of the questions

- 3-a. I am both an A and a B,
- 3-b. I am a not-A or a not-B.

If question 3-a is selected with probability p and 3-b with probability \bar{p} , the resulting estimator $\hat{\pi}_{AB}$ (see Warner [2]) of π_{AB} is unbiased and has variance $V(\hat{\pi}_{AB}) = [\pi_{AB} \bar{\pi}_{AB} + f(p)]/n$. It is interesting to note that neither procedure is uniformly better than the other. In fact, if we let $p_1 = p_2 = p$, (for simplicity, and in accordance with Section 5), then

$$V(\hat{\pi}_{AB})/V(\tilde{\pi}_{AB}) = 1 + f(p)[\pi_A + \pi_B + f(p) - 1]/V(\tilde{\pi}_{AB}),$$

which is less than one when $\pi_A + \pi_B + f(p) < 1$, and this occurs when π_A , π_B and $1/2 - |1/2 - p|$ are all relatively close to zero. Both procedures provide relatively little confidentiality under these circumstances (p is close to zero or one), but the two question approach is revealing for all the responses in AUB, while the one question case is revealing only for those in A and B. Thus, the smaller variance of $\hat{\pi}_{AB}$ is likely to be somewhat offset by a greater likelihood of nonresponse

and/or untruthful responses.

Some comparisons between $V(\hat{\pi}_{AB})$ and $V(\tilde{\pi}_{AB})$ are given in Table 1. Note that either variance can be considerably larger than the other, but $V(\tilde{\pi}_{AB})$ is the larger only when sampling for rare characteristics with a small p . More typically, $V(\hat{\pi}_{AB})$ is considerably larger than $V(\tilde{\pi}_{AB})$. This illustrates that as one goes from a single sensitive question to a two sensitive question analysis, there can be a great increase in the variance of the estimate of the joint probability. Such a comparison will show even greater increases in variances when more than two sensitive questions are asked. This can be easily seen for the special case of independent characteristics.

Table 1. $V(\hat{\pi}_{AB})$ and $V(\tilde{\pi}_{AB})$ for Selected Values of π_A , π_B , π_{AB} , and p

π_A	π_B	π_{AB}	p	$nV(\hat{\pi}_{AB})$	$nV(\tilde{\pi}_{AB})$
.01	.0075	.0025	.4	36.107	6.002
.01	.0075	.0025	.1	0.025	0.143
.04	.0300	.0100	.4	36.430	6.010
.04	.0300	.0100	.1	0.040	0.151
.16	.0400	.0133	.4	37.213	6.013
.16	.0400	.0133	.1	0.061	0.154
.64	.3200	.1067	.4	41.855	6.095

5. Efficiency of the Estimators; Choice of p_1 's

For the two sensitive question case one is usually interested in estimating π_A , π_B and π_{AB} individually. Hence a reasonable measure of the efficiency of the estimation procedure is the trace of the variance-covariance matrix; that is, the quantity

$$v(p_1, p_2) = n[V(\hat{\pi}_A) + V(\hat{\pi}_B) + V(\hat{\pi}_{AB})]. \quad (5.1)$$

It is to be noted that $v(p_1, p_2)$ depends on p_1 and p_2 only through the function $f(p)$. Further, $f(p)$ has the properties: $f(0) = 0$; $f(\bar{p}) = f(p)$, $0 \leq p \leq 1$; $df/dp > 0$, $0 < p < 1/2$; and $f(p) \rightarrow \infty$ as $p \rightarrow 1/2$. It follows that $\min_{p_1, p_2} v(p_1, p_2) = v(0, 0)$, which is the measure obtained for the direct question approach.

Since maximum statistical efficiency cannot be achieved without destroying all confidentiality, one could take the approach of selecting p_1 and p_2 to achieve a given preassigned efficiency. That is, given efficiency $1/r$, $r > 1$, select p_1 and p_2 to satisfy

$$v(p_1, p_2) = rv(0, 0). \quad (5.2)$$

This is one equation in two unknowns and therefore has infinitely many solutions. Without loss of generality assume $0 < p_1 < 1/2$ and $0 < p_2 < 1/2$. Then, solving (5.2) for p_2 as a function of p_1 we obtain

$$p_2 = [1 - (4f_2 + 1)^{-1/2}]/2,$$

where

$$f_2 = [(r-1)v(0, 0) - f_1(1 + \pi_B)]/[f_1 + 1 + \pi_A],$$

and $f_1 = f(p_1)$. Routine application of the calculus establishes that p_2 is a continuous, strictly decreasing and concave function of p_1 .

An obvious and convenient solution for equation (5.2) can be obtained by choosing $p_1 = p_2$ (or $p_1 = \bar{p}_2$, because of the symmetry of $f(p)$). This solution maximizes $\min(p_1, p_2)$. That is, it gives the greatest protection to the respondent on each question for a given efficiency. It was conjectured that there may occasionally be an advantage in choosing $p_1 \neq p_2$, especially when π_A differed greatly from π_B . It was felt that a small decrease in one of the p 's would lead to a far larger increase in the other. This did not happen as can be seen from the entries of Table 2, in which $r = 10$. For example, when $\pi_A = .64$ and $\pi_B = .01$, with $\pi_{AB} = .00125$, then $p_1 = p_2 = .237$ is one solution to (5.2). Another solution is $p_1 = .220$ and $p_2 = .249$. That is, a sacrifice in p of .017 on question 1 yields a gain of .012 on question 2. There is little advantage in terms of protecting the respondent's privacy, in choosing the latter solution over $p_1 = p_2 = .237$. In sum, the solution $p_1 = p_2$ is quite reasonable and little can be gained by choosing $p_1 \neq p_2$.

Let us now consider the solution $p_1 = p_2 = p$ (say). Table 3 gives solutions of $v(p, p) = rv(0, 0)$ for selected values of $1/r$, and π_A , π_B and π_{AB} . Typical choices for p recommended in the literature (see [2] for details) are in the range from .2 to .3. In the two sensitive question case, Table 3 shows that this will typically result in a loss in efficiency of at least 70%, compared to direct questioning and assuming all responses are truthful.

Table 2. Values of p_1 and p_2 that Achieve 10% Efficiency for Selected π_A , π_B and π_{AB}

π_A	.16	.32	.64			
π_B	.16	.08	.01			
π_{AB}	.04	.04	.00125			
	P_1	P_2	P_1	P_2	P_1	P_2
	.000	.346	.000	.342	.000	.301
	.036	.342	.036	.339	.034	.297
	.069	.338	.071	.335	.068	.293
	.104	.333	.107	.329	.101	.288
	.138	.325	.142	.323	.135	.281
	.173	.316	.178	.313	.169	.271
	.208	.302	.213	.300	.203	.258
	.242	.282	.249	.281	.220	.249
	.263	.263	.267	.267	.237	.237
	.277	.249	.284	.249	.271	.206
	.311	.185	.320	.186	.304	.148
	.346	.000	.356	.000	.339	.000

Table 3. Values of p Required for Given Efficiencies for Selected π_A , π_B and π_{AB}

π_A	π_B	π_{AB}	Efficiency ($1/r$)			
			.8	.4	.2	.1
.05	.05	.0125	.012	.061	.122	.187
.10	.05	.0250	.018	.082	.153	.219
.20	.15	.0750	.037	.131	.211	.273
.25	.05	.0375	.027	.112	.190	.255
.25	.25	.0625	.038	.142	.223	.284
.25	.25	.2500	.047	.163	.244	.301
.40	.05	.0250	.029	.118	.197	.262
.55	.25	.1250	.042	.152	.234	.294
.75	.05	.0250	.022	.096	.172	.240
.75	.70	.5250	.041	.150	.234	.295

6. Effects of Untruthful Responses

Randomized response schemes are designed to reduce bias due to nonresponse or lying to protect one's privacy. To investigate the effects of such lying let

$$t_A = P(\text{an A tells the truth about being an A}),$$

and define t_B and t_{AB} similarly. Assume further, that persons not in a sensitive group will not claim to be members of such a group. Then for either randomized response or the direct question approach we have

$$v_A = E(\hat{\pi}_A) = \pi_A t_A \quad (6.1)$$

$$v_B = E(\hat{\pi}_B) = \pi_B t_B \quad (6.2)$$

$$v_{AB} = E(\hat{\pi}_{AB}) = \pi_{AB} t_{AB} \quad (6.3)$$

and the biases are therefore:

$$b(\hat{\pi}_A) = \pi_A \bar{t}_A \quad (6.4)$$

$$b(\hat{\pi}_B) = \pi_B \bar{t}_B \quad (6.5)$$

$$b(\hat{\pi}_{AB}) = \pi_{AB} \bar{t}_{AB} \quad (6.6)$$

It should be noted that the values of the t 's in equations (6.1)-(6.6) will usually differ for the direct question case as compared to randomized response. Hopefully, the t 's will have higher values for the randomized response scheme.

We will compare our randomized response scheme with direct questioning by using the sum of the mean squared errors (MSE):

$$M_t(p_1, p_2) = v(p_1, p_2) + b^2(\hat{\pi}_A) + b^2(\hat{\pi}_B) + b^2(\hat{\pi}_{AB})$$

where $M_t(0, 0)$ is the measure obtained by direct questioning and $v(p_1, p_2)$ is given by equation (5.1).

Tables 4-7 give values of the sum of the biases ("bias") and $M_t(p_1, p_2)$ for selected

values of t and $p_1 = p_2 = p$. Notice that for truthful responses by both methods, the randomized response technique leads to far larger values of MSE than does direct questioning. Consider now the case $\pi_A = .16$, $\pi_B = .12$, $\pi_{AB} = .04$, $p = .3$, $n = 1000$, and truthful responses for the Warner scheme while $t_A = .7$, $t_B = .6$, $t_{AB} = .5$ for direct questioning. Then $M_t(0,0) = .0051939$ while $M_t(.3, .3) = .00499$. So in this extreme case the Warner scheme is superior to direct questioning. If we consider the above example with $t_A = .9$, $t_B = .7$, $t_{AB} = .7$ then $M_t(0,0) = .0019234$ and direct questioning would be superior. In summary, the Warner scheme is superior to direct questioning only when randomized responses would produce considerable increases in the rate of truthful responses. It should be noted that in the above discussion we have excluded consideration of different rates of nonresponse resulting from using the two approaches.

References

- [1] Horvitz, D. G., Bernard G., and Abernathy, J. R., "Recent Developments in Randomized Response Design," in J. N. Srivastava, ed., A Survey of Statistical Design and Linear Models, New York: North-Holland Publishing Company, 1975, 271-85.
- [2] Warner, S. L., "Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias," Journal of the American Statistical Association, 60 (March 1965), 63-9.

Table 4. Effects of Lying on MSE; $n = 100$,
 $\pi_A = .04$, $\pi_B = .01$, $\pi_{AB} = .00667$

t_A	t_B	t_{AB}	bias	$M_t(0,0) \times 10^6$	$M_t(p,p) \times 10^6$
$p = .3 \quad p = .1$					
1.0	1.0	1.0	0.0	549	44682 3630
1.0	0.9	0.8	.00233	529	44648 3608
0.9	0.7	0.7	.00900	492	44532 3563
0.7	0.6	0.5	.01933	536	44458 3594
0.6	0.4	0.2	.02733	608	44451 3658

Table 5. Effects of Lying on MSE; $n = 1000$,
 $\pi_A = .04$, $\pi_B = .01$, $\pi_{AB} = .00667$

t_A	t_B	t_{AB}	bias	$M_t(0,0) \times 10^7$	$M_t(p,p) \times 10^7$
$p = .3 \quad p = .1$					
1.0	1.0	1.0	0.0	549	44682 3630
1.0	0.9	0.8	.00233	554	44674 3633
0.9	0.7	0.7	.00900	753	44794 3824
0.7	0.6	0.5	.01933	2076	45999 5134
0.6	0.4	0.2	.02733	3492	47336 6541

Table 6. Effects of Lying on MSE; $n = 100$,
 $\pi_A = .16$, $\pi_B = .12$, $\pi_{AB} = .04$

t_A	t_B	t_{AB}	bias	$M_t(0,0) \times 10^6$	$M_t(p,p) \times 10^6$
$p = .3 \quad p = .1$					
1.0	1.0	1.0	0.0	2784	49935 6188
1.0	0.9	0.8	.020	2825	49819 6212
0.9	0.7	0.7	.064	3970	50439 7301
0.7	0.6	0.5	.116	6867	52758 10136
0.6	0.4	0.2	.168	11708	57075 14921

Table 7. Effects of Lying on MSE; $n = 1000$,
 $\pi_A = .16$, $\pi_B = .12$, $\pi_{AB} = .04$

t_A	t_B	t_{AB}	bias	$M_t(0,0) \times 10^7$	$M_t(p,p) \times 10^7$
$p = .3 \quad p = .1$					
1.0	1.0	1.0	0.0	2784	49935 6188
1.0	0.9	0.8	.020	4697	51691 8084
0.9	0.7	0.7	.064	19234	65703 22565
0.7	0.6	0.5	.116	51939	97830 55208
0.6	0.4	0.2	.168	104444	149811 107657