

# ADJUSTING FOR NONRESPONSE BIAS OF CORRELATED ITEMS USING LOGISTIC REGRESSION<sup>1</sup>

Pao-Sheng and Robin Fisher  
Robin Fisher, U.S. Bureau of the Census, Washington, D.C. 20233

## 1. Introduction

Item nonresponse in sample surveys is the failure to obtain a specific question that should have been answered. In particular, item nonresponse rate on the survey of the economic is usually high. This can result for many reasons, the most frequent being "refusals to answer", which can relate to the underlying data value (non-ignorable nonresponse) and human behavior. Most designs of the survey questionnaire incorporate procedures for following up on missing responses to items identified as either especially important to the overall quality of the survey data or with previously noted high nonresponse rate. For example, the design of the Survey of Income and Program Participation (SIPP) questionnaire incorporated procedures for following up on missing responses to the items of wage and salary income, income received from self-employment and interest and dividend income. The response status on these items by the same individual are most likely correlated. The problem of missing items for categorical variable has been examined from the perspective of modeling the mechanisms of nonresponse by Fay (1986), Chambers and Welsh (1993), Alho (1990), and Särndal (1981).

This paper proposes a method of adjusting item nonresponse in presence of callback based on a generalized logistic regression model that can account for the correlation among responses on items. The probability of response for any item is represented by a logistic regression model, in which the value of that item, the response status of the rest of the items and the available covariates, which may include the observed item variables for all the individuals by the last callback, are explanatory variables. The respondents are assumed to answer some or all of items after one or more call-backs. The parameters of our model can be estimated by taking a conditional maximum likelihood approach based on the respondents. This approach has the advantage of the simple expression of conditional logistic model. The estimated individual probabilities of responding are used in a Horvitz-Thompson type estimator to reduce

bias in the estimation of sample means for every single item.

## 2. The Logistic Regression for Correlated Responses

### 2.1 A Class of Conditional Logistic Models

Let  $I = \{1, \dots, n\}$  be a set of indices for  $n$  individuals selected in a simple random sample. Let  $X_i = (X_{i1}, \dots, X_{iL})$  be the set of item outcomes from individual  $i$  and they suffer the nonresponse,  $i=1, \dots, n$ , where  $X_{ij}$  expresses the outcome of the  $l^{\text{th}}$  item from individual  $i$ , the value of which becomes known when individual  $i$  responds for the item  $l$ . The vector of covariates of individual  $i$  is denoted by  $Z_i$ . Suppose up to  $J \geq 2$  attempts are made to capture the data for an individual. Define the nonresponse indicator vector  $U_{ij} = (U_{ij1}, \dots, U_{ijL})^T$ , where, for  $l=1, \dots, L$  with  $U_{ijl} = 1$  if individual  $i$  was captured at the  $j^{\text{th}}$  attempt for the  $l^{\text{th}}$  item, and  $U_{ijl} = 0$  otherwise.

Define  $y_{ijl} = \sum_{k=1}^j U_{ikl}$  ( $i=1, \dots, n; l=1, \dots, L$ ) for

short. Then  $y_{ijl} = 1$ , if and only if individual  $i$  was captured by the  $j^{\text{th}}$  attempt for the  $l^{\text{th}}$  item. If  $U_{ijl}$  are correlated, the probability for  $U_{ijl}=1$  not only depends upon  $X_{ij}$  and  $Z_i$  but also depends on the responses for the rest of the items. First, consider the class of conditional logistic models when  $j=1$

$$\begin{aligned} \text{logit Pr } (U_{i1l} = 1 | U_{i1}^{-l}, X_{i1}, Z_i) \\ = F_{1l}(U_{i1}^{-l}) + G_{1l}(X_{i1}, Z_i) \end{aligned} \quad (1)$$

when  $j > 1$

$$\begin{aligned} \text{logit Pr } (U_{ijl} = 1 | U_{ij}^{-l}, S_{ij-1}, X_{i1}, A_i) \\ = F_{jl}(U_{ij}^{-l}) + G_{jl}(X_{i1}, Z_i) \text{ if } y_{ijl} = 0 \end{aligned} \quad (2)$$

$\text{Pr } (U_{ijl} = 1 | U_{ij}^{-l}, S_{ij-1}, X_{i1}, Z_i) = 0$  if  $y_{i,j-1,l} = 1$   
where  $U_{ij}^{-l}$  is  $U_{ij}$  with the exclusion of  $U_{ijl}$ , for  $j=1, \dots, J$ ,  $S_{ij-1} = (U_{i1}^T, \dots, U_{i,j-1}^T)^T$ .

<sup>1</sup> This paper reports the general results of research undertaken by Census Bureau staff. The views expressed are attributable to the authors and do not necessarily reflect those of the Census Bureau.

$F_{il}$  is an arbitrary function of  $U_{il}^{-1}$  such that

$$u_{ijl} = \begin{cases} 0, & \text{if } y_{ij-1,l} = 0 \\ 1 & \text{if } y_{ij-1,l} = 1 \end{cases}$$

$\sum_{l=1}^L U_{il} F_{il} (I_{il}^{-1})$  is invariant under permutation of

$U_{il}$ 's, where  $I_{il}^{-1}$  is  $U_{il}^{-1}$  with  $U_{ilk} = 0$  for

$k > 1$ ,  $F_{jl}$  is an arbitrary function of  $U_{jl}^{-1}$

such that  $\sum_{l=1}^L U_{jl} F_{jl} (I_{jl}^{-1}, S_{i,j-1})$  is invariant under permutation of  $U_{jl}$ 's, and

where  $I_{jl}^{-1}$  is  $U_{jl}^{-1}$  with  $U_{jlk} = 0$  for  $k > 1$

Thus  $F_{jl}$  is a function describing the dependence of item  $l$  on the response status of the other items in and before call-back attempt  $j$ . The function  $G$  describes the dependence on the outcomes  $X$  and the covariates  $Z$ .

For given  $F_{il}$  and  $F_{jl}$ , from (1) and (2) we have the joint probability of  $U_{ij}$ 's uniquely defined as when  $j = 1$

$$\begin{aligned} & Pr(U_{il} | X_p, Z_i) \\ & = \exp \left[ \sum_{l=1}^L U_{il} [F_{il}(I_{il}^{-1}) + G_{il}(X_p, Z_i)] \right] / d_{il} \end{aligned} \quad (3)$$

where

$$d_{il} = \sum_{u_{i1}, \dots, u_{iL}} \exp \left[ \sum_{l=1}^L u_{il} [F_{il}(I_{il}^{-1}) + G_{il}(X_p, Z_i)] \right],$$

$$u_{iil} = 0, 1, l=1, \dots, L$$

$$\begin{aligned} & \text{when } j \geq 1 \\ & Pr(U_{ij} | S_{i,j-1}, X_p, Z_i) \\ & = \exp \left[ \sum_{l=1}^L U_{ijl} [F_{jl}(I_{jl}^{-1}, S_{i,j-1}) + G_{jl}(X_p, Z_i)] \right] / d_{ij} \end{aligned} \quad (4)$$

where

$$d_{ij} = \sum_{u_{ij1}, \dots, u_{ijL}} \exp \left[ \sum_{l=1}^L u_{ijl} [F_{jl}(I_{jl}^{-1}, S_{i,j-1}) + G_{jl}(X_p, Z_i)] \right]$$

Equations (3) and (4) follow from an argument similar to that given Liang and Zeger (1989) in the appendix. Note that  $d_{il}$  and  $d_{ij}$  are the normalizing constants which involve a sum of  $2^L$  exponential terms.

## 2.2 An example

1. A special case is that where the response probability of item  $l$  in attempt  $j$  depends on the responses of the other item only through their number

in attempt  $j$  (denoted by  $r_{ijl} = \sum_{k \neq l} U_{ijk}$  and their

number by attempt  $j-1$  (denoted by

$$t_{ij-1,l} = \sum_{k \neq l} y_{ij-1,k}). \text{ That is, when } j = 1$$

$$\begin{aligned} & \text{logit } Pr(U_{iil} = 1 | U_{ii}^{-1}, X_p, Z_i) \\ & = F_1(r_{iil}) + G_{il}(X_p, Z_i) \end{aligned}$$

when  $j > 1$

$$\begin{aligned} & \text{logit } Pr(U_{ijl} = 1 | U_{ij}^{-1}, S_{i,j-1}, X_p, Z_i) \\ & = F_j(r_{ijl}, t_{ij-1,l}) + G_{jl}(X_p, Z_i) \text{ if } y_{ij-1,l} = 0 \end{aligned}$$

$$Pr(U_{ijl} = 1 | U_{ij}^{-1}, S_{i,j-1}, X_p, Z_i) = 0 \text{ if } y_{ij-1,l} = 1$$

When  $j=1$ , we have the joint probability of  $U_{ij}$  uniquely defined as

$$\begin{aligned} & Pr(U_{il} | X_p, Z_i) \\ & = \exp \left[ \sum_{l=1}^L U_{il} [F_1(B_{iil}) + G_{il}(X_p, Z_i)] \right] / d_{il} \end{aligned} \quad (5)$$

$$\text{where } B_{iil} = \sum_{k=1}^{l-1} U_{iik} \text{ and we assume that } B_{iil} = 0$$

Similarly, when  $j > 1$ , the joint probability of  $U_{ij}$  conditional on  $S_{i,j-1}$  is uniquely defined as

$$G_{jl}(X_{ij}, Z_i) = \alpha_{jl} + X_{ij}\beta_{1j} + Z_i^T\beta_{2j} \quad (7)$$

$$PR(U_{ij} | S_{i,j}, X_{ij}, Z_i) = \exp \left[ \sum_{l=1}^L U_{ijl} [F_j(B_{ijl} t_{i,j-1,l}) + G_{jl}(X_{ij}, Z_i)] \right] / d_{ij}$$

where  $U_{ijl} = 0$  when  $y_{i,j-1,l} = 1$   
 $= 0$  otherwise

(6)

$$B_{ijl} = \sum_{k=1}^{l-1} U_{ijk} \quad \text{and we assume that } B_{ij1} = 0$$

We might choose the model to describe a situation in which some respondents are more willing to respond than others. Consider item  $l$  in attempt  $j$ . If the subject has already responded to a large number of items, we might expect him or her to be more likely to respond to item  $l$ . We index the function  $F$  with the call-back attempt to allow the dependence to change with the call-back number.

This type of model has been considered by Qu, William, Beck and Goormastic (1987) for the case  $J=1$ .

$$\text{If we let } \sum_{t=0}^{T_i} F_l(t) = 0 \quad \text{when } T_{ij} = 0, \text{ and}$$

$$\sum_{t=t_{j-1}}^{T_{ij}} F_j(t) = 0 \quad \text{when } T_{ij} = t_{j-1} \text{ and } F_j(B_{ijl}, t_{i,j-1}) = F_j(B_{ijl} + t_{i,j-1}), l=1, \dots, L,$$

We get the following class of conditional logistic models.

When  $j=1$

$$Pr(U_{i1l}=1 | U_{i1}^{-1}, X_{i1}, Z_i) = \frac{e^{F_1(r_{i1}) + G_{1l}(X_{i1}, Z_i)}}{1 + e^{F_1(r_{i1}) + G_{1l}(X_{i1}, Z_i)}}$$

When  $j > 1$

$$Pr(U_{ijl} = 1 | U_{ij}^{-1}, S_{i,j-1}, X_{ij}, Z_i) = \frac{e^{F_j(r_{ij}^*) + G_{jl}(X_{ij}, Z_i)}}{1 + e^{F_j(r_{ij}^*) + G_{jl}(X_{ij}, Z_i)}} \text{ if } y_{i,j-1,l} = 0$$

$$= 0 \text{ if } y_{i,j-1,l} = 1$$

Where  $r_{ijl}^* = r_{ijl} + t_{i,j-1,l}$

Here is one simple model.

$$F_j(r_{ijl}^*) = r_{ijl} \delta, \text{ and } F_j(r_{ijl}^*) = r_{ijl}^* \delta \text{ for } j > 1 \quad (8)$$

Therefore the characteristics  $X_i$  and covariates  $Z_i$  affect the capture probabilities for the same way for each attempt. The effect of  $X_i$  is felt on  $P_{ijl}$  only through the  $l^{\text{th}}$  characteristics  $X_{ij}$ . Different attempts may have different capture probabilities depending on the  $\alpha_{ij}$ 's. This may reflect varying methods at callback or the possibility that the respondent's probability of response changes after a number of calls. We can imagine, for example, that a respondent may develop some resistance after even a small number of attempts have been made in which case  $\alpha_{ij}$  decreases in  $j$ . Also notice that the number of responses in other items affect each item the same way.

### 3. Estimation Procedure (Conditional Maximum Likelihood Approach)

Without loss of generality we can order the data so that by observation 1 through  $n_i$  are the responded items for the  $i^{\text{th}}$  individual and  $n_i + 1$  through  $L$  are the nonresponded items, we can estimate the probabilities based on the following 'working' conditional likelihood.

$$L^o = \prod_{i \in (I_A \cup I_S)} \prod_{l=1}^{n_i} \prod_{j=1}^J v_{ijl}^{U_{ijl}}$$

where  $I_A$  denotes the set of individuals who answer all the items,  $I_S$  denotes the set of individuals who only answer some of the items, and

$$v_{ijl} = \frac{\mu_{ijl}}{1 - \mu_{i,j-1,l}}$$

where  $\mu_{i,j-1,l} = 1 - \sum_{k=1}^j \mu_{ikl}$ , and where

$$\mu_{i1l} = P_{i1l} = \frac{e^{F_1(U_{i1}^{-1}) + G_{1l}(X_{i1}, Z_i)}}{1 + e^{F_1(U_{i1}^{-1}) + G_{1l}(X_{i1}, Z_i)}}$$

$$\mu_{ijl} = P_{ijl} \prod_{k=1}^{j-1} (1 - P_{ikl}), \quad j=2, \dots, J, \text{ and where}$$

$$P_{ijl} = \frac{F_{jl}(U_{ij}^{-1}, S_{i,j-1}) + G_{jl}(X_{ij}, Z_i)}{1 + F_{jl}(U_{ij}^{-1}, S_{i,j-1}) + G_{jl}(X_{ij}, Z_i)}$$

Maximum 'working' conditional likelihood estimates of the parameters can be found by numerically maximizing the log of this function with respect to the parameters involved. Consider the assumptions of (9) and (10). Does not have a unique maximum. One way to solve this problem is to use the available additional information in conjunction with the likelihood equation. For computational advantage, we use the approach proposed by Alho (1991).

Let  $I_{1l}$  be the set of individuals captured at the first attempt for the  $l^{\text{th}}$  item ( $l=1, \dots, L$ ),  $I_{2l}$  the set of individuals captured at the second attempt for the  $l^{\text{th}}$  item, etc. Let  $I_{j+1,l}$  be the set of individuals that are not captured for the  $l^{\text{th}}$  item at all. Define  $n_{jl} = \text{card}(I_{jl})$ , for  $l=1, \dots, L$  and  $j=1, \dots, J+1$ ; thus  $n = n_{1l} + \dots + n_{J+1,l}$  for  $l=1, \dots, L$ .

Note that for the  $l^{\text{th}}$  item, given  $S_{i,j-1}$  and  $U_{ij}^{-1}$ , we have, for  $j=1, \dots, J$

$$E \left[ \sum_{i \in \bigcup_{j=1}^J I_{ij}} U_{ij} \frac{1-P_{ij}}{P_{ij}} \mid S_{i,j-1}, U_{ij}^{-1} \right] = (n_{j1} + \dots + n_{j+1,l}) - \sum_{i \in \bigcup_{j=1}^J I_{ij}} P_{ij}$$

Notice  $E(n_{j1} \mid S_{i,j-1}, U_{ij}^{-1}) = \left[ \sum_{i \in \bigcup_{j=1}^J I_{ij}} P_{ij} \right]$  ;

suppose we use  $n_{j1}$  to estimate this. We estimate the expectation on the left-hand side with the observed value also. This yields the likelihood equations

$$\sum_{i \in I_{ij}} \exp(-\alpha_{ji} - W_{ii}^T \beta_l - r_{ij} \delta) = n - (n_{1l} + \dots + n_{jl}), \quad l=1, \dots, L$$

Given  $\beta_l$  and  $\delta$  we can thus solve for  $\alpha_{jl}$  by taking

$$\alpha_{jl} = -\log(n - n_{1l} - \dots - n_{jl}) / \sum_{i \in I_{ij}} \exp(-W_{ii}^T \beta_l - r_{ij} \delta) \quad (j=1, \dots, J, l=1, \dots, L)$$

To solve for

$$\alpha = (\alpha_1^T, \dots, \alpha_J^T)^T \text{ where } \alpha_l = (\alpha_{1l}, \dots, \alpha_{Jl})^T, \text{ for } l = 1, \dots, L$$

$$\beta = (\beta_1^T, \dots, \beta_J^T)^T \text{ and } \delta$$

we use an iteration based on Newton's method.

Differentiating the log likelihood  $L$  with respect to  $\beta_l$ , we get

$$\frac{\partial L}{\partial \beta_l} = - \sum_{j=1}^J \left[ \sum_{i=1}^J (U_{ij} - v_{ij}) \sum_{k=1}^J P_{ik} \right] W_{ii} = 0$$

We can solve numerically for  $\alpha$ ,  $\beta$ , and  $\delta$ .

Having calculated the estimate  $\hat{\alpha}$ ,  $\hat{\beta}$ ,  $\hat{\delta}$ , the Horvitz-Thompson type estimator was considered based on the requirement of unbiasedness. Define  $X_l = (X_{1l}, \dots, X_{nl})^T$ ,  $l=1, \dots, L$ .

The true sample mean for the item  $l$  is

$$\bar{X}_l = X_l^T \mathbf{1}_n / n \text{ where } \mathbf{1}_n \text{ is a vector of } n \text{ ones. By}$$

translating  $\hat{\alpha}$ ,  $\hat{\beta}$ ,  $\hat{\delta}$ , into (7) and (8), we can

calculate the estimates  $\hat{\mu}_{i,j+1,l}$  to get the conditional

unbiased Horvitz-Thompson type estimator of  $\bar{X}_l$  ( $l=1, \dots, L$ ) as

$$\bar{X}_l = \frac{1}{n} \sum_{i \in I_{jL}} X_{il} Y_{iL} / \hat{\mu}_{i,j+1,l}$$

Let  $\hat{\gamma} = (\hat{\alpha}, \hat{\beta}, \hat{\delta})$  be the estimates of

$\gamma = (\alpha, \beta, \delta)$  subject to regularity conditions on the  $X_{il}$ 's and  $Z_i$ 's, the constraints used to ensure identifiability, coverage to continuous relation between the  $\alpha_{ij}$ 's and  $(\beta_p, \delta)$  which is satisfied by the true

parameter vector  $\gamma$ . The proof for the consistency of the parameter estimator  $\hat{\gamma}$  can be given by emulating the standard argument from Fahrmeir & Kaufmann (1985).

#### 4. References

Alho, J. (1990). "Adjusting for Nonresponse Bias Using Logistic Regression." *Biometrika* 77, 617-624.

Besag, J. (1974). "Spatial Interaction and the Statistical Analyses of Lattice Systems." *Journal of the Royal Statistical Society, Ser. B*, 36, 192-236.

Chambers, R. L. and Welsh A. H. (1993). "Log-linear Models for Survey Data with Non-ignorable non-response." *J. R. Statis. Soc. B*, 55, 157-170.

Fay, R. E. (1986). "Causal Model for Patterns of Nonresponse." *J. Am. Statist. Ass.*, 81, 354-365.

Fahrmeir, L and Kaufman, H. (1985). "Consistency and Asymptotic Normality of the Maximum Likelihood Estimator in Generalized Linear Models." *Ann. Statist.* 13, 342-368.

Liang, K. and Zeger, S.L. (1989). "A Class of Logistic Regression Models for Multivariate Binary Time Series." *Journal of the Am. Statist. Assoc.* 84, 447-451.

Qu, Y., Williams, G. W., Beck G. J., & Goormastic, M. (1987). "A Generalized Model of Logistic Regression for Clustered Data." *Commun. Statist. - Theory Meth.* 16, 3447-3476.

Särndal, C.E. (1981). "Frameworks for Inference in Survey Sampling with Applications to Small Area Estimation and Adjustment for Nonresponse." *Bull. Int. Statist. Inst.* 49, 494-513

Zhao, L.P. & Prentice, R.L. (1990). "Correlated Binary Regression Using a Quadratic Exponential Model." *Biometrika* 77, 642-648.