

CROSS-CLASSIFIED DATA FROM COMPLEX SAMPLE SURVEYS:
THE EFFECTS OF COMPLEX ADJUSTMENT PROCEDURES ON VARIANCES

Daniel H. Freeman, Jr.* and Gary G. Koch
University of North Carolina at Chapel Hill

1. Introduction

The United States Bureau of the Census response error model (hereafter, CRE) has been found to be useful for the elucidation of non-sampling errors by a number of investigators (Hansen *et. al.* (1961), Hansen *et. al.* (1964), and Koch *et. al.* (1975)). The CRE extends in a straightforward manner to complex cross-classifications of survey data (Freeman (1975)). This paper will utilize the CRE model to investigate the effects of complex post-sampling data adjustments on the estimated variances.

The major portion of the discussion is devoted to the effect of "raking the data" to known margins. This procedure, proposed by Deming and Stephan (1940), is also known as iterative proportional fitting (denoted IPF). IPF is also used to obtain maximum likelihood estimates for fitting log-linear models to data from a Poisson or multinomial distribution (Fienberg (1970)). An approximation to the variance matrix generated by the procedure is given in Section 3 along with a statistic for the evaluation of the variance reduction achieved by the procedure. These results are then applied to the special case of simple ratio adjustments for representativeness and non-response. These simple ratio adjustments have been investigated empirically by Brock *et. al.* (1975).

While IPF is a post-sampling adjustment aimed at variance reduction, it is also of interest to fit other types of functions to survey data. Examples of these are functions which are directed at the underlying process which generated the responses of the surveyed populations. The multiple logit and the bivariate Weibull are two examples of such functions and are discussed in Section 5. The important point that is brought out is that these functions are within the scope of the CRE and the method of analysis is completely consistent with that used for the IPF discussion.

2. The IPF Post-sampling Adjustment Procedure

One approach to the minimization of non-sampling errors is to adjust the survey estimates of the frequency counts in cross-classified data to correspond to an independently determined set of population parameters of margins. The IPF adjustment uses these margins directly in the adjustment process. Thus the data set is forced to reflect pre-determined criteria of representativeness, an intuitively appealing property. However, in the context of simple random sampling IPF has several additional properties. Fienberg (1970) showed that IPF will converge to a set of domain or cell estimates under fairly general conditions. These conditions are that an underlying multinomial is a correct model and that

*now at Yale University

there are observations in all of the fitted cells. The latter condition is true for many large surveys as long as the cross-classification is not too complex. The importance of the former condition is now known. A related property is that underlying measures of association are preserved by the procedure (Mosteller (1968)). A third property, to be exploited here, is that IPF maximized the likelihood equation of the multinomial distribution (Fienberg (1970)).

Some notation is required for the simplicity of discussion. Let the population consist of N individuals, denoted i , n of which are selected for the sample. The sampling procedure may be characterized by the indicator function, u_i :

$$u_i = \begin{cases} 1 & \text{if individual } i \text{ is selected,} \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

For simplicity it is assumed that there are no complex measurement effects. Implicitly there is a known structure on the expected response vector, Y_i . This means that the category of response, j , reflects the joint response components for \tilde{Y}_i which may be arrayed as an s -dimensional contingency table. The k -th variable may contain up to J_k levels, $k = 1, \dots, s$.

For example, consider a set of 8 binary responses. Then \tilde{Y}_i is a vector with 2^8 components and,

$$Y_{i\tilde{j}} = Y_i(j_1, \dots, j_8) = \begin{cases} 1 & \text{if individual } i \text{ is classified} \\ & \text{in category } (j_1, \dots, j_8) \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Here, $j_k = 1, 2$ for $k = 1, 2, \dots, 8$.

The corresponding estimates for the population distribution are given by,

$$\hat{\tilde{N}} = \sum_{i=1}^N \frac{u_i}{\phi_i} Y_i \quad \text{and} \quad (3)$$

$$\hat{\tilde{N}}' = (\hat{N}_{11\dots 1}, \hat{N}_{11\dots 2}, \hat{N}_{22\dots 1}, \hat{N}_{22\dots 2})$$

where $E\{u_i\} = \phi_i$ is the selection probability of the i -th individual and,

$$\hat{N}_{j_1, j_2, \dots, j_8} = \sum_{i=1}^N \frac{u_i}{\phi_o} Y_i(j_1, j_2, \dots, j_8). \quad (4)$$

The components of various margins may also be expressed as vectors. For example,

$$\hat{N}_{+ \dots j_k \dots j_k \dots +} = \sum_{i=1}^N \frac{u_i}{\phi_i} \sum \dots \sum Y_i(j_1, j_2, \dots, j_8) \quad (5)$$

where the sum is over all j with $j_k' = j_k''$, $j_k'' = j_k''$. Then

$$\hat{N}_{+...k'...k''...+} = (\hat{N}_{+...1...1...+}, \hat{N}_{+...1...2...+}, \dots, \hat{N}_{+...j_{k'}...j_{k''}...+}).$$

For the general situation suppose the $(s-2, s-1)$, $(s-2, s)$ and $(s-1, s)$ pairwise margins may be taken as given. Then the IPF adjustment is given by,

$$\begin{aligned} & N_{j_1, j_2, \dots, j_s}^{*(I, 1)} \\ &= N_{j_1, j_2, \dots, j_s}^{*(I-1, 3)} \frac{N_{+...+, j_{s-2}, j_{s-1}, +}}{N_{+...+, j_{s-2}, j_{s-1}, +}^{*(I-1, 3)}}; \\ & N_{j_1, j_2, \dots, j_s}^{*(I, 2)} \\ &= N_{j_1, j_2, \dots, j_s}^{*(I, 1)} \frac{N_{+...+, j_{s-2}, +, j_s}}{N_{+...+, j_{s-2}, +, j_s}^{*(I, 1)}}; \\ & N_{j_1, j_2, \dots, j_s}^{*(I, 3)} \\ &= N_{j_1, j_2, \dots, j_s}^{*(I, 2)} \frac{N_{+...+, j_{s-1}, j_s}}{N_{+...+, j_{s-1}, j_s}^{*(I, 2)}} \end{aligned} \quad (6)$$

where $N_{j_1, j_2, \dots, j_s}^{*(0, 3)} = \hat{N}_{j_1, j_2, \dots, j_s}$ and I is the superscript indicating the cycle of the adjustment. This procedure converges in most cases of interest to a set of estimates \tilde{N}^* , where

$$\begin{aligned} \tilde{N}_{+...s-2, s-1, +}^* &= \tilde{N}_{+...s-2, s-1, +}^*, \\ \tilde{N}_{+...s-2, +, s}^* &= \tilde{N}_{+...s-2, +, s-1}^*, \\ \tilde{N}_{+...+, s-1, s}^* &= \tilde{N}_{+...+, s-1, s}^*, \end{aligned} \quad (7)$$

the fixed margins. These estimates maximize the multinomial likelihood equation

$$\begin{aligned} & L(\tilde{N} | \tilde{N}_{+...s-2, s-1, +}, \tilde{N}_{+...s-2, +, s-1}, \\ & \quad \tilde{N}_{+...+, s-1, s}) = \ln(\text{constant}) \\ & + \sum_{\text{all } j} \sum_{\text{all } j} \hat{N}_{j_1, j_2, \dots, j_s} \ln \{ \pi_{j_1, j_2, \dots, j_s} \} \\ & - N \left(\sum_{\text{all } j} \pi_{j_1, j_2, \dots, j_s} - 1 \right), \end{aligned} \quad (8)$$

where $\pi_{j_1, j_2, \dots, j_s} = N_{j_1, j_2, \dots, j_s} / N$ is to be estimated subject to (7). Thus, a set of estimates \tilde{P}^* are obtained from IPF, $\pi \hat{=} \tilde{P}^* = \frac{1}{N} \tilde{N}^*$, where " $\hat{=}$ " means is estimated by.

There are two reasons for being interested in the characterization of \tilde{P}^* given by (8). The first (Koch and Tolley (1975)) is that they correspond to a set of smoothed means. The corresponding equations (8) are a useful modeling tool which are computationally simple, intuitively suggestive, and have desirable large sample properties. The second reason is that (8) may well have meaning in the context of super population theory and that the only difficulty with \tilde{N}^* is in obtaining estimates of the variance matrix which reflect the complex sampling plan characterized by the u_i . As will be shown, if consistent estimates of the variance matrix $\hat{\tilde{N}}$ exist, then corresponding estimates of \tilde{N}^* exist. This point follows easily from the theory of implicit functions.

The constraints in (7) imply an underlying set of parameters for the multiplicative associations in the cross-classification (Fienberg (1970)). That is to say, there exists an implicit function $F(\beta)$ where the β are a unknown set of parameters, such that

$$\begin{aligned} \pi &= F(\beta), \quad \beta' = (\beta_1, \beta_2, \dots, \beta_u), \\ u &\leq \left(\prod_{k=1}^S J_k \right) - 1. \end{aligned} \quad (9)$$

Thus, $F(\beta)$ may be substituted for π in (8) and (8) maximized under (7) by IPF. Thus, IPF solves,

$$\left[\frac{\partial}{\partial \beta} (\ln F(\beta)) \right]' \left[F(\beta) - \hat{p} \right] = 0 \quad (10)$$

where $\hat{p} = \frac{1}{N} \hat{\tilde{N}}$. It follows that β is implicitly defined in terms of \hat{p} , $\hat{\beta} = H(\hat{p})$. For regular functions $F(\beta)$, $H(\hat{p})$ is differentiable about some value of π , say π_0 , and therefore can be expanded in a Taylor series about π_0 ,

$$\hat{\beta} = H(\hat{p}) \hat{=} h(\hat{p}) = h(\pi_0) + \left[\frac{\partial h(\hat{p})}{\partial \hat{p}} \right]_{\hat{p}=\pi_0} (\hat{p} - \pi_0). \quad (11)$$

From (3) it follows that $E\{\hat{p}\} = \pi_0$, so $h(\hat{p})$ is a consistent estimate of β when π_0 obtains and the IPF model is appropriate. This permits the definition of the linearized variance of $\hat{\beta}$ to be,

$$\tilde{V}(\hat{\beta}) = \left[\frac{\partial \tilde{H}(\hat{p})}{\partial \hat{p}} \right]_{\hat{p}=\pi_0} \tilde{V}(\hat{p}) \left[\frac{\partial \tilde{H}(\hat{p})}{\partial \hat{p}} \right]_{\hat{p}=\pi_0}^T \quad (12)$$

where $\tilde{V}(\hat{p})$ may be any valid and consistent estimate of the covariance matrix of \hat{p} which

reflects the underlying sample design. Such estimates may be obtained by the method of Balanced Repeated Replication (BRR) as discussed by McCarthy (1966), Koch *et. al.* (1975), and Brock *et. al.* (1975). Once $\tilde{V}(\hat{\beta})$ is in hand another similar Taylor series expansion of $F(\beta)$ yields $\tilde{V}(\hat{p}^*)$ since (9) implies $\hat{p}^* = F(\hat{\beta})$. Thus $\tilde{V}(\hat{p}^*)$ is sufficient for a consistent estimate of $\tilde{V}(\hat{p}^*)$ as was to be shown. The only difficult term is $(\partial \tilde{H}(\hat{p})/\partial \hat{p})$ which can be found by differentiating (10) by \hat{p} .

3. The Variance Structure of the IPF Model

The equation for $\tilde{V}(\hat{\beta})$ shows the importance of the form which $F(\beta)$ takes. Specifically, for the model shown in (8) which IPF maximized, $F(\beta)$ takes a log-linear form,

$$\pi = \exp\{X\beta\}, \quad (13)$$

where X is a $((\prod_{j,k=1}^s J_k) \times u)$ design matrix which

corresponds to a complete factorial design (Fienberg (1970b)), with additive constraints. For example, if there are three binomial responses, then

$$\tilde{X} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & -1 & 1 & -1 & -1 & -1 \\ 1 & 1 & -1 & 1 & -1 & 1 & -1 & -1 \\ 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 \\ 1 & -1 & 1 & 1 & -1 & -1 & 1 & -1 \\ 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 \\ 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 \\ 1 & -1 & -1 & -1 & 1 & 1 & 1 & -1 \end{bmatrix}, \text{ and } \tilde{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \\ \beta_7 \\ \beta_8 \end{bmatrix} \quad (14)$$

If β_1 through β_8 are estimated, this is equivalent to not adjusting the data. If the total population size is known then β_1 is fixed. If the first order margins are known then β_2, β_3 , and β_4 are fixed. If the second order margins are known then β_5, β_6 , and β_7 are fixed. If the third order margin is known then the sample need not have been drawn. These parameters correspond to various levels of multiplicative association. The usual level of adjustment is done in the

reverse of the order of model reduction in the linear models context. That is, the highest order terms are the last to be eliminated.

Returning to the $\tilde{V}(\hat{\beta})$ implied by (13) the solution for $[\partial \tilde{H}(\hat{p})/\partial \hat{p}]$ follows the method of proof given in Koch and Tolley (1975).

Let,

$$F_{i,j}^{(1)}(\hat{\beta}) = \left[\frac{\partial}{\partial \beta_j} (F_i(\beta)) \right]_{\beta=\hat{\beta}} \quad (15)$$

where $i = 1, \dots, r = \prod_{k=1}^s J_k$, $j = 1, \dots, u$, and

$$F_{i,j}^{(2)}(\hat{\beta}) = \left[\frac{\partial^2}{\partial \beta_j^2} (F_i(\beta)) \right]_{\beta=\hat{\beta}} \quad (16)$$

Noting that $\sum_{i=1}^r F_i(\beta) = 1$ implies $\sum_{i=1}^r F_{i,j}^{(1)}(\beta) = 0$

$= \sum_{i=1}^r F_{i,j}^{(2)}(\beta) = 0$. Then

$$\left[\frac{\partial \tilde{H}(\hat{p})}{\partial \hat{p}_i} \right]_{\hat{p}=\pi_1} = \frac{F_{i,j}^{(1)}(\hat{\beta})/F_i(\hat{\beta})}{\sum_{i=1}^r [F_{i,j}^{(1)}(\hat{\beta})]^2/F_i(\hat{\beta})} \quad (17)$$

for $i = 1, \dots, r$; $j = 1, \dots, u$.

Substituting (15) and (16) into (13) yields:

$$\begin{aligned} F_{i,j}^{(1)}(\hat{\beta})/F_i(\hat{\beta}) &= x_{ij}, \\ [F_{i,j}^{(1)}(\hat{\beta})]^2/F_i(\hat{\beta}) &= x_{ij}^2 p_j^*. \end{aligned} \quad (18)$$

These in turn yield

$$\begin{aligned} & \left[\frac{\partial \tilde{H}(\hat{p})}{\partial \hat{p}} \right]_{\hat{p}=\hat{p}^*} \\ &= \begin{bmatrix} \frac{x_{11}}{\sum_{i=1}^r x_{i1}^2 p_i^*} & \frac{x_{21}}{\sum_{i=1}^r x_{i1}^2 p_i^*} & \dots & \frac{x_{r1}}{\sum_{i=1}^r x_{i1}^2 p_i^*} \\ \frac{x_{12}}{\sum_{i=1}^r x_{i2}^2 p_i^*} & \frac{x_{22}}{\sum_{i=1}^r x_{i2}^2 p_i^*} & \dots & \frac{x_{r2}}{\sum_{i=1}^r x_{i2}^2 p_i^*} \\ \dots & \dots & \dots & \dots \\ \frac{x_{1u}}{\sum_{i=1}^r x_{iu}^2 p_i^*} & \frac{x_{2u}}{\sum_{i=1}^r x_{iu}^2 p_i^*} & \dots & \frac{x_{ru}}{\sum_{i=1}^r x_{iu}^2 p_i^*} \end{bmatrix} = G(\hat{p}^*) \end{aligned} \quad (19)$$

When $J_k = 2$, $k = 1, \dots, s$ then $x_{ij}^2 = 1$ all i, j so

$$\left[\frac{\partial H(\hat{p})}{\partial \hat{p}} \right] \bigg|_{\hat{p}=\hat{p}^*} = \tilde{X}' \quad (20)$$

So returning to the original problem of approximating $V(\hat{p}^*)$:

$$\hat{\beta} = H(\hat{p}) , \quad (21)$$

$$V(\hat{\beta}) = G(\hat{p}^*) V(\hat{p}) [G(\hat{p}^*)]' , \quad (22)$$

$$\hat{p}^* = \exp\{X \hat{\beta}\} , \quad (23)$$

and

$$V(\hat{p}^*) = D_{\hat{p}}^* G' G V(\hat{p}) G' G D_{\hat{p}}^* , \quad (24)$$

where $G = G(\hat{p}^*)$ is defined in (19) and $D_{\hat{p}}^*$ is a diagonal matrix with the IPF estimates on the main diagonal. As noted in Section 2, $V(\hat{p})$ may be obtained from BRR or any other method which yields consistent estimates.

Given these estimates of the variance matrix of \hat{p}^* it is desirable to find an overall measure of the variance reduction achieved by IPF. One approach is suggested by the use of weighted least squares as discussed in Koch *et. al.* (1975).

Let \hat{p}_1 be the vector of sample estimates with the r -th element deleted and $V_1(\hat{p}_1)$ be the corresponding covariance matrix. Similarly define \hat{p}_u^* and $V_u(\hat{p}_u^*)$ as the vector of adjusted estimates with the last $r - u$ elements deleted and the corresponding covariance matrix. Then if $\hat{p}_1 = (\hat{p}_1 - 1 b_1)$ where b_1 is estimated by weighted least squares and 1 is a vector of $(r - 1)$ 1's, then

$$Q_1 = \hat{p}_1' V_1^{-1}(\hat{p}_1) \hat{p}_1 \quad (25)$$

is a measure of the total variation in the unadjusted cross-classification. Similarly,

$$Q_u = \hat{p}_u' V_u^{-1}(\hat{p}_u^*) \hat{p}_u \quad (26)$$

is a measure of the total variation in the adjusted cross-classification. For large data sets in the super-population context Q_1 and Q_u have χ^2 distributions with $r - 1$ and $u - 1$ degrees of freedom respectively. Moreover, these form valid indices of the overall variance reduction achieved by IPF.

4. Application to Poststratification

Some sample surveys are conducted in such a way that for single margins involving several variables complete information about the target population is available. For example, in a survey to determine health characteristics, the demographic characteristics for each of the survey strata may be independently determined.

Thus, the ratio of the known variables to their estimated values may be used to inflate the strictly estimated survey variables. This process is known as poststratification. Frequently, this inflation takes the form of a simple multiplicative adjustment.

The poststratification procedure is in fact a simplified form of the IPF procedure discussed in Sections 2 and 3. What is done is to array the demographic variables as a single variable with a number of levels equal to the product of the number of levels in each of the known variables. Thus, the poststratification corresponds to a single iteration of the IPF algorithm. It immediately follows from the discussion in Section 3 that the Q statistic is increased to the extent of the elimination of sampling variance and measurement error. There is a corresponding reduction in the standard errors. A cautionary note should be observed in that poststratification, as with the IPF procedure, will induce bias in the estimate to the extent of error in the known variables.

The important point is that when poststratification, or for that matter IPF, is contemplated the Q statistics and their corresponding ratios may be computed to estimate the relative improvement in precision. Thus, assessments may be made and weighed against the possibility of inducing bias in the estimation process.

5. Complex Postsampling Models

Subsequent to the selection and analysis of a sample survey, substantive analysts often become aware of the possibility that the responses were in fact generated by an underlying stochastic process. These processes usually can be formalized into some type of statistical function. Two such functions are the multiple logit and the Weibull distributions. Each of these distributions have found wide acceptance in the biological and social sciences because of their flexibility and interpretability. The fitting of these distributions for survey data may easily be implemented in the framework of the CRE model as will be demonstrated. Here the multiple logit is considered in the context of multiple binomials and the Weibull is considered in the context of a bivariate distribution for a cross-sectional representative sample.

Consider the situation where the survey responses are classified by two levels for each of three factors. This is again a 2^3 design on the response space. Let π_{ijk} denote the probability of a response at joint levels (i, j, k) where $i = 1, 2$, $j = 1, 2$, $k = 1, 2$; and

$$\sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 \pi_{ijk} = 1 . \quad (27)$$

Then the multiple logit model is:

$$\ln \pi_{ijk} = x_{ijk}^{(1)} \beta_1 + x_{ijk}^{(2)} \beta_2 + x_{ijk}^{(3)} \beta_3 + x_{ijk}^{(4)} \beta_4 + x_{ijk}^{(5)} \beta_5 + x_{ijk}^{(6)} \beta_6 - \ln D \quad (28)$$

where

$$x_{ijk}^{(1)} = 1 \quad k = 2, 0 \quad \text{otherwise};$$

$$x_{ijk}^{(2)} = 1 \quad j = 2, 0 \quad \text{otherwise};$$

$$x_{ijk}^{(3)} = 1 \quad i = 2, 0 \quad \text{otherwise};$$

$$x_{ijk}^{(4)} = x_{ijk}^{(2)} x_{ijk}^{(3)}$$

$$x_{ijk}^{(5)} = x_{ijk}^{(2)} x_{ijk}^{(4)}$$

$$x_{ijk}^{(6)} = x_{ijk}^{(3)} x_{ijk}^{(4)}$$

and

$$D = \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 \exp \left\{ \sum_{r=1}^6 x_{ijk}^{(r)} \beta_r \right\}.$$

Implicit in the model is the constraint (27) and the eighth degree of freedom is reserved for a goodness of fit test. The formulation of Sections 2 and 3 in terms of models which correspond to maximum likelihood functions may also be used. The essential point is the use of maximum likelihood search procedures to estimate the β_r .

The estimation of standard errors again requires the evaluation of

$$\frac{\partial}{\partial \beta} \{ \ln [F(\beta)] \} \quad (29)$$

where,

$$\ln F_{ijk}(\beta) = \sum_{r=1}^6 x_{ijk}^{(r)} \beta_r - \ln D. \quad (30)$$

In matrix notation the solutions may be written as in (10).

The covariance matrix of the parameters may now be obtained from (17) and (12). In particular, after some algebra

$$\frac{\partial H(p)}{\partial p} \bigg|_{p=\pi} = \begin{bmatrix} \frac{-1}{\pi_{..1}} & \frac{1}{\pi_{..2}} & \frac{-1}{\pi_{..1}} & \frac{1}{\pi_{..2}} & \frac{-1}{\pi_{..1}} & \frac{1}{\pi_{..2}} & \frac{-1}{\pi_{..1}} & \frac{1}{\pi_{..2}} \\ \frac{-1}{\pi_{.1.}} & \frac{-1}{\pi_{.1.}} & \frac{1}{\pi_{.2.}} & \frac{1}{\pi_{.2.}} & \frac{-1}{\pi_{.1.}} & \frac{-1}{\pi_{.1.}} & \frac{1}{\pi_{.2.}} & \frac{1}{\pi_{.2.}} \\ \frac{-1}{\pi_{1..}} & \frac{-1}{\pi_{1..}} & \frac{-1}{\pi_{1..}} & \frac{-1}{\pi_{1..}} & \frac{1}{\pi_{2..}} & \frac{1}{\pi_{2..}} & \frac{1}{\pi_{2..}} & \frac{1}{\pi_{2..}} \\ \frac{-1}{1-\pi_{.22}} & \frac{-1}{1-\pi_{.22}} & \frac{-1}{1-\pi_{.22}} & \frac{1}{\pi_{.22}} & \frac{-1}{1-\pi_{.22}} & \frac{-1}{1-\pi_{.22}} & \frac{-1}{1-\pi_{.22}} & \frac{1}{\pi_{.22}} \\ \frac{-1}{1-\pi_{2.2}} & \frac{-1}{1-\pi_{2.2}} & \frac{-1}{1-\pi_{2.2}} & \frac{-1}{1-\pi_{2.2}} & \frac{1}{\pi_{2.2}} & \frac{1}{\pi_{2.2}} & \frac{-1}{1-\pi_{2.2}} & \frac{1}{\pi_{2.2}} \\ \frac{-1}{1-\pi_{22.}} & \frac{-1}{1-\pi_{22.}} & \frac{-1}{1-\pi_{22.}} & \frac{-1}{1-\pi_{22.}} & \frac{-1}{1-\pi_{22.}} & \frac{-1}{1-\pi_{22.}} & \frac{1}{\pi_{22.}} & \frac{1}{\pi_{22.}} \end{bmatrix} \quad (31)$$

Thus, as with the IPF procedure, $\hat{\pi}$ and $V(\hat{\pi})$ may be obtained by substituting \hat{p} for π in (28). This in turn yields the estimates of p^* and $V(p^*)$ under the model, where p^* refers to the logit model proportions.

An asymptotically equivalent procedure is that of weighted least squares. The formulation is as follows. Let

$$K \ln A \pi \hat{=} X \beta \quad (32)$$

where

$$A = I_8, \quad K = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{bmatrix}$$

Then β may be estimated directly from weighted least squares procedures and corresponding estimates of π obtained by inverting (32). Again, the only critical aspect is in obtaining $V(\hat{p})$ from some procedure such as BRR. Thus, it is apparent that functions with a simple exponential form, such as multiple logits, may be fitted to complex survey data. The bivariate Weibull distribution represents a more complex function.

The Weibull distribution has been of interest to statisticians since its introduction to the literature, partly because of its usefulness in data which are thought to reflect either increasing or decreasing hazard functions. If t represents the time to the occurrence of an event of interest (e.g., a death or the detection of a tumor), then the Weibull cumulative distribution function may be written as:

$$G(t|\mu, \delta, w) = 1 - \exp[-\mu(t-w)^\delta] \quad (33)$$

for $t \geq w$ and $\mu, \delta \geq 0$.

As demonstrated in Freeman, Freeman, and Koch (1974) such models may easily be fitted in the usual weighted least squares framework. This may be seen by examining $\theta(t)$ which denotes the log-log of the survival rate. Specifically let

$$\theta(t) = \ln[-\ln\{1 - G(t|\mu, \delta, w)\}] = \ln(\mu) + \delta \ln(t-w), \quad (34)$$

where $\ln(\mu)$ and δ are to be estimated and w is to be fixed by the experimental situation. With this reparameterization $\theta(t)$ can easily be fitted, both for univariate and bivariate data.

Typically, the Weibull distribution is fitted to data which is assumed to arise from a simple random sample. However, the CRE model and the variance estimation procedure of BRR permit the fitting of this distribution to data arising from a complex sample survey. The only change would be instead of estimating $\text{Var}(\hat{p})$ based on the multinomial model assumption, $V(\hat{p})$ would be based on the Balanced Repeated Replication procedure.

6. Summary and Conclusions

This paper was concerned with various post-sampling adjustments to complex survey data and the fitting of complex functions to such data. The IPF (iterative proportional fitting) procedure was investigated and an approximation to the variance structure induced by the procedure was developed. Further, a statistic, Q , was suggested for evaluating the overall variance reduction achieved by the procedure. This statistic is particularly useful since it is a natural by-product of the inference structure discussed in Koch *et. al.* (1975).

The discussion of IPF was followed by a discussion of the procedure for fitting more complex distributions to data. The multiple logit and bivariate Weibull distributions were examined in this context. It was pointed out that once the distribution of interest had been reparameterized into the linear model framework, the only change from the simple random sampling inference structure was in terms of estimating the variance matrix. This was shown to be easily accomplished when BRR is used. Thus, it follows that the response error model in conjunction with the BRR permits the implementation of the weighted least squares methodology for inference in the post-sampling analysis of complex sample survey data.

ACKNOWLEDGMENTS

This research was in part supported by the National Institutes of Health (Grants GM-70004 and HD-00371) and by the U. S. Bureau of the Census through Joint Statistical Agreements JSA 74-2 and JSA 75-2.

REFERENCES

- Brock, D.B., Freeman, D.H., Freeman, J.L. and Koch, G.G. (1975). An application of categorical data analysis to the National Health Interview Survey, Proceedings of the 1975 Social Statistics Section of the ASA.
- Deming, W.E. and Stephan, F.F. (1940). On a least squares adjustment of a sampled frequency table when the expected marginal totals are known, Ann. Math. Stat. **11**, 427-444.
- Fienberg, S.E. (1970). An iterative procedure for estimation in contingency tables, Ann. Math. Stat. **41**, 907-918.
- Freeman, D.H. (1975). The regression analysis of data from complex sample surveys: an empirical investigation of covariance matrix estimation, Institute of Stat. Mimeo Series No. 1020, Chapel Hill, N. C.
- Freeman, D.H., Freeman, J.L., and Koch, G.G. (1974). A modified χ^2 approach for fitting Weibull models to synthetic life tables, Institute of Stat. Mimeo Series No. 958, Chapel Hill, N. C.
- Hansen, M.H., Hurwitz, W.N. and Pritzker, L. (1964). The estimation and interpretation of gross differences and the simple response variance, in Contributions to Statistics, Pergamon Press Ltd., London, 111-136.
- Koch, G.G., Freeman, D.H. and Freeman, J.L. (1975). Strategies in the multivariate analysis of data from complex surveys, Int. Stat. Rev. **43**, 55-74.
- Koch, G.G. And Tolley, H.D. (1975). A generalized modified χ^2 analysis of categorical bacterial survival data from a complex dilution experiment, Biometrics **31**, 59-92.
- McCarthy, P.J. (1966). Replication: an approach to the analysis of data from complex surveys, Vital and Health Statistics, P.H.S. Pub. No. 1000, Ser. 2, No. 14, National Center for Health Statistics, Rockville, Md.
- Mosteller, F. (1968). Association and estimation in contingency tables, J.A.S.A. **63**, 1-28.