

VARIANCE ESTIMATION FROM COMPLEX SURVEY DATA USING LINERAZATION METHOD: IMPACT OF DAVID BINDER

A. Demnati¹ and J.N.K. Rao²

ABSTRACT

David Binder's pioneering 1983 paper published in the *International Statistical Review* provided a unified linearization approach to variance estimation from complex survey data. This important paper stimulated much new research on linearization variance estimation in subsequent years. This talk will review some of that work, including the alternative linearization method of Demnati and Rao (2004, 2010). We apply our approach to the estimation of total variance of calibration estimators of the generalized linear models parameters when missing items have been imputed using either deterministic or random imputation.

KEY WORDS: Best linear predictor; Finite population parameters; Imputation; Item nonresponse; Model parameters; Multiple weight adjustments.

RÉSUMÉ

L'article novateur de David Binder publié en 1983 dans le *International Statistical Review* a présenté une approche unifiée de l'estimation de la variance par linéarisation à partir de données d'enquête complexes. Cet article important a fait souffler un vent nouveau sur la recherche sur l'estimation de la variance par linéarisation dans les années qui ont suivi. Nous passerons en revue certains de ces travaux, y compris la méthode de recharge par la linéarisation de Demnati et Rao (2004, 2010). Nous appliquons notre approche à l'estimation de la variance totale des estimateurs par calage pour les paramètres des modèles linéaires généralisées lorsque les réponses manquantes ont été imputées de façon déterministe ou aléatoire.

MOTS CLÉS: Meilleur prédicteur linéaire; paramètres de population finie; imputation; non-réponse partielle; paramètres du modèle; ajustements multiple de poids.

1. INTRODUCTION

Traditionally, statistical agencies collect data to estimate relatively simple finite population quantities such as totals or means. Estimators and associated variance estimators for such parameters are well studied in the literature. Currently, survey data are also used to investigate relationships between variables using statistical models, such as linear and logistic regression models, and interest is in the estimation of the model parameters and associated variance estimators.

Using Taylor series linearization approach, Binder (1983) presented a unified method of estimating the design-based variance of estimators of "census" model parameters based on complex sample designs from finite populations. This method is

¹ A. Demnati, Business Survey Methods Division, Statistics Canada, Ottawa, Canada, Abdellatif.Demnati@statcan.gc.ca

² J.N.K. Rao, School of Mathematics and Statistics, Carleton University, Ottawa, Canada, JRao@math.carleton.ca

particularly useful when the parameters are defined implicitly as solutions to estimating equations. Taylor linearization is generally applicable to any sampling design that permits unbiased variance estimation for linear estimators. However, it can lead to multiple variance estimators that are asymptotically design unbiased under repeated sampling. The choice among the variance estimators, therefore, requires other considerations such as (i) approximate unbiasedness for the model variance of the estimator under an assumed model, (ii) validity under a conditional repeated sampling framework. For example, in the context of simple random sampling and the ratio estimator, $\hat{Y}_R = (\bar{y}/\bar{x})X$, of the population total Y , Royall and Cumberland (1981) showed that a commonly used linearization variance estimator, $\mathcal{G}_{LI} = N^2(n^{-1} - N^{-1})s_z^2$, does not track the conditional variance of \hat{Y}_R given \bar{x} , unlike the jackknife variance estimator \mathcal{G}_J . Here \bar{y} and \bar{x} are the sample means, X is the known population total of the auxiliary variable x , s_z^2 is the sample variance of the residuals $z_k = y_k - (\bar{y}/\bar{x})x_k$ and n denotes the sample size. By linearizing the jackknife variance estimator, \mathcal{G}_J , a different linearization variance estimator, $\mathcal{G}_{JL} = (\bar{X}/\bar{x})^2 \mathcal{G}_{LI}$, is obtained. This variance estimator also tracks the conditional variance as well as the unconditional variance, where $\bar{X} = X/N$ is the mean of x . As a result, \mathcal{G}_{JL} or \mathcal{G}_J may be preferred over \mathcal{G}_{LI} . Särndal, Swensson and Wretman (1989) showed that \mathcal{G}_{JL} is both asymptotically design unbiased and asymptotically model unbiased in the sense of $E_m(\mathcal{G}_{JL}) \approx \text{Var}_m(\hat{Y}_R)$, where $\text{Var}_m(\hat{Y}_R)$ is the model variance of \hat{Y}_R under a “ratio model”: $E_m(y_k) = \beta x_k$; $k = 1, \dots, N$ and the y_k 's are independent with model variance $\text{Var}_m(y_k) = \sigma^2 x_k$, $\sigma^2 > 0$. Thus, \mathcal{G}_{JL} is a good choice from either the design-based or the model-based perspective. Binder (1996) presented an elegant “cookbook” approach to Taylor linearization that leads directly to \mathcal{G}_{JL} -type linearization variance estimators. He applied the method to smooth functions of estimated totals, $g(\hat{Y}_1, \dots, \hat{Y}_m)$, generalized regression estimators and the Wilcoxon rank sum statistic.

Demnati and Rao (2004) proposed an alternative Taylor linearization approach to variance estimation that is theoretically justifiable and at the same time leads directly to a \mathcal{G}_{JL} -type variance estimator for general designs. They applied the method under the design based approach to a variety of problems, covering regression calibration estimators of a total Y and other estimators defined either explicitly or implicitly as solutions of estimating equations. They obtained a new variance estimator for a general class of calibration estimators that includes generalized raking ratio and generalized regression estimators. They also extended the method to two-phase sampling and obtained a sampling variance estimator that makes fuller use of the first phase sample data compared to traditional linearization variance estimators. Demnati and Rao (2010) studied total variance estimation in the context of finite populations assumed to be generated from superpopulation models and analytical inferences on model parameters are of interest. If the sampling fractions are negligible, then the sampling variance captures almost the entire variation generated by the design and model random processes. However, when the sampling fraction is not negligible, the model variance should be taken into account in order to construct valid inferences on model parameters under both randomization processes.

Demnati and Rao (2002) extended their method to the scalar case of missing responses when weight adjustments for complete nonresponse and imputation for item nonresponse based on smooth functions of observed values, in particular ratio imputation, are used. A standard approach for handling multivariate missing items is to treat one variable at a time and

perform separate imputation. Although separate imputations make the data complete, it does not take advantage of correlations between observed items. Suppose that the respondent is unable to provide a value for each item but rather the respondent is willing to provide the total of two or more items. For example, in business surveys large enterprises are unable to provide value for each combination of industry and geography. However, the total revenue or expense at business level is readily available. In this situation, the separate imputation problem is further complicated by the availability of totals that have to be satisfied by the items. The total value of numerical items must be equal to the sum of its parts. Additionally, if only one item has a missing value then the imputation method should compute the correct value of this variable from the observed values of the other items and their total. This is an example of a situation where the imputation procedure should determine the missing value uniquely from the observed values.

In section 2, we briefly review the Binder method for variance estimation, while in section 3, we give a brief account of the Demnati-Rao (DR) method for total variance estimation. In section 4, we consider estimating equations when calibration and imputation for item nonresponse have been used, and study estimators obtained as solutions to estimating equations. Finally in section 5, we give a preview of our recent work on multivariate imputation for general patterns of missingness under observed control totals. Our multivariate imputation method preserves automatically the observed items while satisfying observed control totals.

2. BINDER'S PIONEERING APPROACH TO VARIANCE ESTIMATION

Consider a finite population of N units identified by a set of indices $P = \{1, \dots, k, \dots, N\}$. Usually, many variables (say M variables of interest) are under study for a given survey. Let $\mathbf{y}_k = (y_{1k}, \dots, y_{Mk})^T$ be the vector of values of the variable of interest $\mathbf{y} = (y_1, \dots, y_M)^T$ attached to unit k , and let $(\mathbf{x}_k^T, \mathbf{t}_k^T, \mathbf{I}_k^T, \mathbf{z}_k^T, \dots)^T$ be the vector of auxiliary variables attached to unit k , where the superscript T denotes the transpose of a vector or a matrix. We denote a typical variable of interest by y , and a typical auxiliary variable by x . Suppose that the model mean of the response y_k is specified by $E_m(y_k) = \mu_k(\mathbf{x}_k^T \boldsymbol{\beta})$, where $\mathbf{x}_k = (x_{1k}, \dots, x_{pk})^T$ is a $p \times 1$ vector of explanatory variables, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ is the $p \times 1$ vector of model parameter and E_m denotes model expectation. Binder (1983) considered the finite population parameter $\boldsymbol{\beta}_N$ defined as solution to ‘‘census’’ estimating equation of the form

$$\mathbf{S}_\beta(\boldsymbol{\beta}; \mathbf{I}_N, \mathbf{y}, \mathbf{A}_x) = \sum_k s_\beta(\boldsymbol{\beta}; y_k, \mathbf{x}_k) - \mathbf{v}_\beta(\boldsymbol{\beta}) = \mathbf{0}, \quad (2.1)$$

where \mathbf{I}_N is the N -vector of 1's, \mathbf{A}_x is a $p \times N$ matrix, with k^{th} column \mathbf{x}_k , $\mathbf{y} = (y_1, \dots, y_N)^T$, \sum_k denote the sum over all the population units, $s_\beta(\boldsymbol{\beta}; y_k, \mathbf{x}_k)$ is a p -dimensional vector-valued function of y_k and \mathbf{x}_k , and the known function $\mathbf{v}_\beta(\boldsymbol{\beta})$ allows for explicitly defined parameters. For linear and logistic regression models, $s_\beta(\boldsymbol{\beta}; y_k, \mathbf{x}_k) = \mathbf{x}_k(y_k - \mu_k(\mathbf{x}_k^T \boldsymbol{\beta}))$, and $\mathbf{v}_\beta(\boldsymbol{\beta}) = \mathbf{0}$. For the special case of the finite population total $Y = \sum_k y_k$, $s_\beta(\boldsymbol{\beta}; y_k, \mathbf{x}_k) = y_k$, $\mathbf{v}_\beta(\boldsymbol{\beta}) = \beta_N$ and $\beta_N = Y$. An estimator $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}_N$ is defined by the design weighted estimating equation

$$\mathbf{S}_\beta(\boldsymbol{\beta}; \mathbf{d}(s), \mathbf{y}, \mathbf{A}_x) = \sum_k d_k(s) s_\beta(\boldsymbol{\beta}; y_k, \mathbf{x}_k) - \mathbf{v}_\beta(\boldsymbol{\beta}) = \mathbf{0}, \quad (2.2)$$

where $d_k(s) = a_k(s) / \pi_k$ denotes the sampling design weights attached to unit k , $a_k(s)$ is the sample membership indicator variable for unit k , $\pi_k = E_p(a_k(s))$ is the inclusion probability for unit k , E_p denotes expectation with respect to the sampling design, and $\mathbf{d}(s) = (d_1(s), \dots, d_N(s))^T$.

Using a Taylor series linearization argument based on the first order approximation, Binder's sampling variance of $\hat{\beta}$ is given by

$$\text{Var}_p(\hat{\beta}) = [\mathbf{J}(\beta_N)]^{-1} \Sigma_v(\beta_N) [\mathbf{J}(\beta_N)]^T, \quad (2.3)$$

where $\mathbf{J}(\beta) = -\partial \mathbf{S}_\beta(\beta; \mathbf{I}_N, \mathbf{y}, \mathbf{A}_x) / \partial \beta$, and $\Sigma_v(\beta)$ is the sampling variance of $\sum_k d_k(s) s(\beta; y_k, \mathbf{x}_k)$. Binder's estimator of $\text{Var}_p(\hat{\beta})$ is

$$\hat{\text{Var}}_p(\hat{\beta}) = [\hat{\mathbf{J}}(\hat{\beta})]^{-1} \hat{\Sigma}_v(\hat{\beta}) [\hat{\mathbf{J}}(\hat{\beta})]^T, \quad (2.4)$$

where $\hat{\Sigma}_v(\hat{\beta})$ is a consistent estimator of $\Sigma_v(\beta_N)$ and $\hat{\mathbf{J}}(\beta) = -\partial \mathbf{S}_\beta(\beta; \mathbf{d}(s), \mathbf{y}, \mathbf{A}_x) / \partial \beta$. Binder (1983) gave regularity conditions for the validity of (2.3) and (2.4), and also showed how the results can be applied to some generalized linear models, including the logistic regression model.

There are a number of competing estimators of the variance that are asymptotically equivalent to (2.4). Binder (1996) presented a "cookbook" approach which produces one of the most favored of those variance estimators. Consider the estimator $\hat{\beta}$ of the finite population parameter β_N which can be expressed as a smooth function

$$\hat{\beta} = g(\hat{Y}_1, \dots, \hat{Y}_M), \quad (2.5)$$

of estimated totals $\hat{\mathbf{Y}} = (\hat{Y}_1, \dots, \hat{Y}_M)^T$, where $\hat{Y}_i = \sum_k d_k(s) y_{ik}$ is an estimator of the population total $Y_i = \sum_k y_{ik}$, $i = 1, \dots, M$, and $\beta_N = g(Y_1, \dots, Y_M)$. Binder' cookbook approach may be summarized as follows:

- 1) Take the total differential $d\hat{\beta} = \sum_i \frac{dg(\hat{Y}_1, \dots, \hat{Y}_M)}{d\hat{Y}_i} d\hat{Y}_i$ of $\hat{\beta}$.
- 2) Replace the total differential $d\hat{\beta}$ by deviations of estimators from their respective population parameters, e.g., $d\hat{\beta}$ is changed to $\hat{\beta} - \beta_N$, and $d\hat{Y}_i$ is changed to $\hat{Y}_i - Y_i = \sum_k d_k(s) y_{ik} - Y_i$; and so on, to get

$$\hat{\beta} - \beta_N = \sum_i \frac{dg(\hat{Y}_1, \dots, \hat{Y}_M)}{d\hat{Y}_i} (\sum_k d_k(s) y_{ik} - Y_i). \quad (2.6)$$

- 3) Rewrite the result of step 2 given by (2.6) as $\hat{\beta} \approx \sum_k z_k d_k(s) + R$, where R does not depends explicitly on the design weight $d_k(s)$.
- 4) Finally obtain the estimated variance as $\mathcal{G}(\hat{\beta}) \approx \hat{\text{Var}}_p(\sum_k d_k(s) z_k) = \hat{\text{Var}}_p(\hat{Z})$, using the formula for the variance estimator of the estimated total \hat{Z} under the specified sampling design, where $\hat{Z} = \sum_k d_k(s) z_k$.

The main difference of this formulation from the standard Taylor linearization method is that in the latter approach the partial derivatives are evaluated at their expected values before z_k is derived, and the unknown parameters in the resulting z_k are then replaced by their estimators.

3. DEMNATI-RAO LINERAZITAION METHOD

3.1 The Method

To motivate the DR method, we start with the case of $\hat{\beta}$ defined as a smooth function (2.5) of estimated totals. DR express $\hat{\beta}$ and β_N as $\hat{\beta} = f(\mathbf{d}(s), \mathbf{A}_y)$ and $\beta_N = f(\mathbf{I}_N, \mathbf{A}_y)$, where \mathbf{A}_y is a $M \times N$ matrix, with k^{th} column $\mathbf{y}_k = (y_{1k}, \dots, y_{Mk})^T$, $k = 1, \dots, N$, and $\mathbf{I}_N = E_p(\mathbf{d}(s))$. Note that $\hat{\beta}$ is a function of both $\mathbf{d}(s)$ and \mathbf{A}_y , but we drop \mathbf{A}_y for simplicity and write $\hat{\beta} = f(\mathbf{d}(s))$. Taylor linearization of $\hat{\beta}$ around \mathbf{Y} gives the approximation

$$\hat{\beta} - \beta_N \approx \partial g(\hat{\mathbf{Y}}) / \partial \hat{\mathbf{Y}}^T |_{\hat{\mathbf{Y}}=\mathbf{Y}} (\hat{\mathbf{Y}} - \mathbf{Y}) = \sum_k \tilde{z}_k (d_k(s) - 1), \quad (3.1)$$

where $\tilde{z}_k = f(\mathbf{b}) / \partial b_k |_{\mathbf{b}=\mathbf{I}_N}$ and \mathbf{b} is a $N \times 1$ vector of arbitrary real numbers. Note that $(\partial g(\hat{\mathbf{Y}}) / \partial \hat{\mathbf{Y}}^T)(\hat{\mathbf{Y}} - \mathbf{Y})$ is equals to the right side of (2.6). The DR sampling variance estimator of $\hat{\beta}$ is given by

$$\mathcal{G}_{DR}(\hat{\beta}) = \mathcal{G}(\sum_k d_k(s) z_k), \quad (3.2)$$

with

$$z_k = \partial f(\mathbf{b}) / \partial b_k |_{\mathbf{b}=\mathbf{d}(s)}. \quad (3.3)$$

Demnati and Rao (2004) gave justification of (3.1), (3.2) and (3.3). Note that DR do not first evaluate the partial derivatives of $\partial f(\mathbf{b}) / \partial b_k$ at $\mathbf{b} = \mathbf{I}_N$ and then substitute estimates of the unknown components. The DR method, therefore, is similar in spirit to Binder's cookbook approach.

Estimation of the finite population parameters, $\beta_N = \mathbf{h}(\mathbf{A}_y)$, or model parameters, $\beta_M = \mathbf{h}(\beta)$, under an assumed super-population model on y are often of interest, where β is the $p \times 1$ super-population model parameter. We now consider a general formulation of the Demnati and Rao (2004, 2010) approach to deriving Taylor linearization variance estimators. This formulation will cover both finite population (or census) parameters, β_N , and model parameters, β_M .

Based on a sample, an estimator, $\hat{\beta}$, is used to estimate both parameters β_N and β_M . Under complete response, the estimator $\hat{\beta}$, obtained as the solution of the estimating equation given by (2.2) is often used as an estimator of the finite population parameter β_N defined as the solution to (2.1) and also the model parameter, β_M , under an assumed model on y .

Let $\mathbf{d}_k = (\mathbf{d}_{1k}^T, \mathbf{d}_{2k}^T, \dots, \mathbf{d}_{gk}^T)^T$ be a $G \times 1$ vector of random weights and \mathbf{u}_k be a $G \times p$ vector of constants for $k = 1, \dots, N$. Let $\hat{\mathbf{U}} = \sum_k \mathbf{u}_k \mathbf{d}_k$ be a linear estimator and, using an operator notation, let $\mathcal{G}(\mathbf{u})$ denote the estimator of variance of $\hat{\mathbf{U}}$. We write $\hat{\beta}$ as $f(\mathbf{A}_d)$, where \mathbf{A}_d is a $G \times N$ matrix with k^{th} column \mathbf{d}_k . The DR linearization variance estimator of $\hat{\beta} = f(\mathbf{A}_d)$ is simply given by $\mathcal{G}_{DR}(\hat{\beta}) = \mathcal{G}(\mathbf{z})$, where $\mathcal{G}(\mathbf{z})$ is obtained from $\mathcal{G}(\mathbf{u})$ by replacing \mathbf{u}_k by $\mathbf{z}_k = \partial f(\mathbf{A}_b) / \partial \mathbf{b}_k |_{\mathbf{A}_b=\mathbf{A}_d}$, where \mathbf{A}_b is a $G \times N$ matrix of arbitrary real numbers with k^{th} column $\mathbf{b}_k = (b_{1k}, \dots, b_{Gk})^T$. The choice of \mathbf{A}_d depends on the random processes involved. Suppose first that the parameter of interest is β_N and we use the estimator given by (2.2). In this case, $g = G = 1$,

$\mathbf{d}_k = d_k = d_k(s)$, and $g(\mathbf{u})$ is the estimator of $\text{Var}(\sum_k \mathbf{u}_k d_k(s))$ with $\mathbf{z}_k = [\hat{\mathbf{J}}(\hat{\boldsymbol{\beta}})]^{-1} s_{\beta}(\hat{\boldsymbol{\beta}}; y_k, \mathbf{x}_k)$, Suppose on the other hand, we are interested in the model parameter $\boldsymbol{\beta}_M$. In this case, $g = 1$, $G = p$, $\mathbf{d}_k = d_k(s) s_{\beta}(\boldsymbol{\beta}; y_k, \mathbf{x}_k)$, and $g(\mathbf{u})$ is the estimator of $\text{Var}(\sum_k \mathbf{u}_k d_k)$ with $\mathbf{z}_k = [\hat{\mathbf{J}}(\hat{\boldsymbol{\beta}})]^{-1}$.

The variance estimators associated with the finite population parameter $\boldsymbol{\beta}_N$ and the model parameter $\boldsymbol{\beta}_M$ are different. In the former case, we estimate the variance by an estimator of the design variance $\text{Var}_p(\hat{\boldsymbol{\beta}}) = E_p(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_N)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_N)^T$ of $\hat{\boldsymbol{\beta}}$, while in the later case, we estimate the total variance

$$\text{Var}(\hat{\boldsymbol{\beta}}) = E_m \text{Var}_p(\hat{\boldsymbol{\beta}}) + \text{Var}_m E_p(\hat{\boldsymbol{\beta}}),$$

where Var_m denotes the variance with respect to the model. It remains to derive $\mathbf{z}_k = \partial f(\mathbf{A}_b) / \partial \mathbf{b}_k |_{\mathbf{A}_b = \mathbf{A}_d}$, when $\hat{\boldsymbol{\beta}}$ is the solution of the estimating equation given by (2.2). One may use a fast approach to derive \mathbf{z}_k , which consists of two steps: a) first derive $\mathbf{z}_{p;k} = \partial f(\mathbf{b}) / \partial \mathbf{b}_k |_{\mathbf{b} = \mathbf{d}(s)}$ where $\hat{\boldsymbol{\beta}} = f(\mathbf{d}(s))$; then b) isolate each component z_{ik} of \mathbf{z}_k corresponding to component d_{ik} of \mathbf{d}_k , after approximating $\text{Var}(\hat{\boldsymbol{\beta}})$ by $\text{Var}(\sum_k \mathbf{z}_{p;k} d_k(s))$. Taking the derivative, we get

$$\mathbf{z}_{p;k} = \partial f(\mathbf{b}) / \partial \mathbf{b}_k |_{\mathbf{b} = \mathbf{d}(s)} = [\hat{\mathbf{J}}(\hat{\boldsymbol{\beta}})]^{-1} s_{\beta}(\hat{\boldsymbol{\beta}}; y_k, \mathbf{x}_k). \quad (3.4)$$

It follows from (3.4) that, when the parameter of interest is the finite population parameter $\boldsymbol{\beta}_N$, the linearized variable is given by (3.4), while in the case where the parameter of interest is the model parameter $\boldsymbol{\beta}_M$, the linearized variable is given $\mathbf{z}_k = [\hat{\mathbf{J}}(\hat{\boldsymbol{\beta}})]^{-1}$.

3.2 Calibration Estimator

Calibration estimator, $\tilde{\boldsymbol{\beta}}$, is the solution to weighted estimating equations of the form

$$\mathbf{S}_{\beta}(\boldsymbol{\beta}; \mathbf{w}(s), \mathbf{y}, \mathbf{A}_x) = \mathbf{0},$$

with weights

$$w_k(s) = d_k(s) F(\mathbf{t}_k^T \hat{\boldsymbol{\gamma}}),$$

and satisfying the calibration constrains

$$\mathbf{S}_{\gamma}(\boldsymbol{\gamma}; \mathbf{d}(s), \mathbf{t}) = \sum_k d_k(s) s_{\gamma}(\boldsymbol{\gamma}; t_k) - \mathbf{v}_{\gamma}(\boldsymbol{\gamma}) = \mathbf{0},$$

with $s_{\gamma}(\boldsymbol{\gamma}; t_k) = F(\mathbf{t}_k^T \hat{\boldsymbol{\gamma}}) t_k$ and $\mathbf{v}_{\gamma}(\boldsymbol{\gamma}) = \mathbf{T}$, where $\mathbf{w}(s) = (w_1(s), \dots, w_N(s))^T$, $\mathbf{t}_k = (t_{1k}, \dots, t_{qk})^T$ and $\mathbf{T} = (T_1, \dots, T_q)^T$ is the vector of known totals of auxiliary variable $\mathbf{t} = (t_1, \dots, t_q)^T$. Calibration estimators are widely used in practice. For example, the choice $F(a) = 1 + a$ gives the generalized regression (GREG) weights and $F(a) = \exp(a)$ leads to raking ratio weights. Taking the derivative $\partial f(\mathbf{b}) / \partial \mathbf{b}_k |_{\mathbf{b} = \mathbf{d}(s)}$ with $\tilde{\boldsymbol{\beta}} = f(\mathbf{d}(s))$, we get

$$\mathbf{z}_{p;k} = [\hat{\mathbf{J}}(\tilde{\boldsymbol{\beta}})]^{-1} F(\mathbf{t}_k^T \hat{\boldsymbol{\gamma}}) \{ s_{\beta}(\tilde{\boldsymbol{\beta}}; y_k, \mathbf{x}_k) - \hat{\mathbf{B}}^T(s) \mathbf{t}_k \}, \quad (3.5)$$

where

$$\hat{\mathbf{B}}(s) = [\sum_k d_k(s) f(\mathbf{t}_k^T \hat{\boldsymbol{\gamma}}) \mathbf{t}_k \mathbf{t}_k^T]^{-1} \sum_k d_k(s) f(\mathbf{t}_k^T \hat{\boldsymbol{\gamma}}) \mathbf{t}_k s_{\beta}^T(\tilde{\boldsymbol{\beta}}; y_k, \mathbf{x}_k),$$

and $f(a) = \partial F(a) / \partial a$.

When the parameter of interest is the finite population parameter β_N , $g = G = 1$, $d_k = d_k(s)$, and the linearized variable is given by (3.5), while in the case where the parameter of interest is the model parameter β_M , $g = 2$, $G = 1 + p$, $d_{1k} = d_k(s)$, $d_{2k} = d_k(s)s_{\beta}(\beta; y_k, \mathbf{x}_k)$, and the components of the linearized variable $\mathbf{z}_k = (z_{1k}^T, z_{2k}^T)^T$ are given by $z_{1k} = -[\hat{\mathbf{J}}(\tilde{\beta})]^{-1}F(\mathbf{t}_k^T \hat{\gamma})\hat{\mathbf{B}}^T(s)\mathbf{t}_k$ and $z_{2k} = [\hat{\mathbf{J}}(\tilde{\beta})]^{-1}F(\mathbf{t}_k^T \hat{\gamma})$.

4. DEMNATI-RAO LINEARIZATION METHOD UNDER IMPUTATION

Item nonresponse occurs in a survey when a sampled unit fails to provide responses to one or more of the survey items. A standard approach for handling missing items is to impute (i.e., fill in) an estimated value for each missing item using an imputation model in combination with an estimate, $\tilde{\beta}$, of the model parameter β , calculated from the incomplete data. We introduce $y_{(k)}^{(o)}$ as the set of observed items of the variables of interest for unit k , and $I_{(k)}^{(o)}$ as the set of all observed values attached to unit k related to both the variables of interest and to the auxiliary variables. In case of complete response from unit k $y_{(k)}^{(o)} = \{y_{1k}, \dots, y_{Mk}\}$, while under complete nonresponse $y_{(k)}^{(o)} = \emptyset$, where \emptyset denotes the empty set. After imputation, it is common practice to treat the imputed values as if they were observed and then compute estimator $\tilde{\theta}$ of the parameter of interest θ as in the complete response case, where θ denotes either a finite population parameter θ_N or a model parameter θ_M . Then it remains to assess the accuracy of the estimator $\tilde{\theta}$ in estimating the parameter θ .

4.1 Conditional Mean Imputation Based Estimator

Let $\tilde{y}_k^{(e)} = E_m(y_k | I_{(k)}^{(o)}; \tilde{\beta})$ be the conditional mean imputed value under the imputation model. Note that $\tilde{y}_k^{(e)} = y_k$ when the item y_k is observed. The previous chapter covers the problem of estimating the variance of $\tilde{\beta}$ solution of the weighted estimating function $S_{\beta}(\beta; \mathbf{w}(s), \mathbf{y}, \mathbf{A}_x) = \mathbf{0}$ having zero mean for the k^{th} component at the true model parameter β , i.e., $E_m\{s_{\beta}(\beta; y_k, \mathbf{x}_k)\} = \mathbf{0}$. In this section estimating equations are used in the more general concept than estimating function which includes weighted log-likelihood estimating function as well as weighted least square estimating functions. A simple example of estimating equation is when one is interested in the overall mean, i.e., $S_{\theta}(\theta; \mathbf{w}(s), \tilde{\mathbf{y}}^{(e)}, \mathbf{A}_x) = \mathbf{0}$ with $s_{\theta}(\theta; \tilde{y}_k^{(e)}, \mathbf{x}_k) = \tilde{y}_k^{(e)} - \theta$ and $\mathbf{v}_{\theta}(\theta) = \mathbf{0}$, when $E_m(y_k) = \mu_k(\mathbf{x}_k^T \beta)$, where \mathbf{x}_k is the vector of auxiliary variables associated with $\tilde{y}_k^{(e)}$ and θ . The parameter of interest induced by the estimator $\tilde{\theta}$ is either the finite population parameter solution to $E_r E_p \{S_{\theta}(\theta; \mathbf{w}(s), \tilde{\mathbf{y}}^{(e)}, \mathbf{A}_x) = \mathbf{0}\}$ or the model parameter solution to $E_m E_r E_p \{S_{\theta}(\theta; \mathbf{w}(s), \tilde{\mathbf{y}}^{(e)}, \mathbf{A}_x) = \mathbf{0}\}$, where E_r denotes expectation with respect to the response mechanism. Given that interest is in estimating the variance of $\tilde{\theta}$, the vector parameter β used for imputation can be seen as a nuisance having $E\{s_{\beta}(\beta; y_k, \mathbf{x}_k)\} = \mathbf{0}$, where $\tilde{\beta}$ is the incomplete-data based estimator solution to

$$S_{\beta}(\beta; \hat{\mathbf{w}}(s), \mathbf{y}, \mathbf{A}_x) = \mathbf{0},$$

with

$$\hat{w}_k(s) = w_k(s)(\delta_k / \tilde{\xi}_k),$$

and $\tilde{\alpha}$ in $\tilde{\xi}_k = \xi_k(\tilde{\alpha})$ is the solution to the logistic regression estimating equation

$$S_{\alpha}(\alpha; \mathbf{w}(s), \boldsymbol{\delta}, \mathbf{A}_I) = \mathbf{0},$$

where $\boldsymbol{\delta} = (\delta_1, \dots, \delta_N)^T$, δ_k is the response indicator for y_k , i.e., $\delta_k = 1$ if y_k is observed and $\delta_k = 0$ if y_k is missing, $s_{\alpha}(\alpha; \boldsymbol{\delta}_k, \mathbf{I}_k) = \mathbf{I}_k(\delta_k - \xi_k)$, $\mathbf{v}_{\alpha}(\alpha) = \mathbf{0}$, $\xi_k = \Pr(\delta_k = 1) = \text{logit}(\mathbf{I}_k^T \alpha)$, and \mathbf{I}_k is the vector of explanatory variables,

After imputation, the estimator $\tilde{\boldsymbol{\theta}}$ is obtained as the solution to

$$S_{\theta}(\boldsymbol{\theta}; \mathbf{w}(s), \tilde{\mathbf{y}}^{(*)}, \mathbf{A}_{\chi}) = \mathbf{0}.$$

The imputed value $\tilde{y}_k^{(*)} = E_m(y_k | \mathbf{I}_{\{k\}}^{(o)}, \tilde{\boldsymbol{\beta}})$ can be rewritten as

$$\tilde{y}_k^{(*)} = \delta_k y_k + (1 - \delta_k) \tilde{y}_k^{(p)},$$

with the missing values replaced by predicted values given by

$$\tilde{y}_k^{(p)} = E_m(y_k | \mathbf{I}_{\{k\}}^{(o)} \setminus \{y_k\}; \tilde{\boldsymbol{\beta}}),$$

available for all sampled units. Note that $\tilde{y}_k^{(*)} = \tilde{y}_k^{(p)}$ for non respondents and $\tilde{y}_k^{(*)} = y_k$ for respondents.

4.2 Random Imputation Based Estimator

Suppose we can sample random values $\tilde{y}_k^{(R)}$ from the conditional distribution with mean $\tilde{y}_k^{(*)}$ and variance $\text{Var}_m(y_k | \mathbf{I}_{\{k\}}^{(o)}, \tilde{\boldsymbol{\beta}})$, so that $\tilde{y}_k^{(R)}$ are independent observations from the known conditional distribution $f(\tilde{y}_k^{(*)}, \text{Var}_m(y_k | \mathbf{I}_{\{k\}}^{(o)}, \tilde{\boldsymbol{\beta}}))$. Under random imputation, the estimator $\tilde{\boldsymbol{\theta}}^{(R)}$ is obtained as the solution to

$$S_{\theta}(\boldsymbol{\theta}; \mathbf{w}(s), \tilde{\mathbf{y}}^{(R)}, \mathbf{A}_{\chi}) = \mathbf{0}.$$

We may write $\tilde{y}_k^{(R)}$ as $\tilde{y}_k^{(R)} = \tilde{y}_k^{(*)} + \varepsilon_k$, where $\varepsilon_k = \tilde{y}_k^{(R)} - \tilde{y}_k^{(*)}$ with $E_R(\varepsilon_k) = 0$, and E_R denotes expectation with respect to the random draw for imputation. In case of conditional mean imputation $\varepsilon_k = 0$.

5. PREVIEW ON MULTIVARIATE CONDITIONAL IMPUTATION

We now give a preview of the recent work of Demnati and Rao (2014) on variance estimation under multivariate conditional imputation for general patterns of missingness under observed control totals. Our multivariate imputation method preserves automatically the observed items while satisfying observed control totals.

5.1 Missing Multivariate Items Situation

Table 1 presents a simple example of a missing data situation for the first six units of a population of size $N = 500$. Three variables and their total, labeled y_1 , y_2 , y_3 and t , are observed, but some of the y_1 , y_2 , y_3 and the t values marked blank in shaded cells are missing. The complete population is generated from a 3-variate Normal distribution with mean $\boldsymbol{\mu} = (5, 10, 15)^T$, variance $\boldsymbol{\sigma}^2 = (5, 10, 15)^T$ and coefficient of correlation $\boldsymbol{\rho} = (\rho_{12}, \rho_{13}, \rho_{23})^T = (.5, .7, .9)^T$. The probabilities of response are set to .5 for

the four variables. The first unit of Table 1 provided complete response, while no information was received from unit 6. Note that the missing item from unit 2 can be determined exactly from the reported values.

Table 1: Observed values

Unit	y_1	y_2	y_3	t
1	7.12	19.71	34.13	60.96
2		-13.53	-7.74	-19.01
3	5.83			30.77
4				19.50
5	10.04			
6				
⋮				

5.2. Observed Data

Consider the case where the data vector of interest $\mathbf{y}_k = (y_{1k}, \dots, y_{Mk})^T$ for unit k is not observed but instead we observe a $m_k \times 1$ vector $\mathbf{I}_k^{(o)}$ which is a linear combination of \mathbf{y}_k ,

$$\mathbf{I}_k^{(o)} = \mathbf{L}_k^T \mathbf{y}_k,$$

where \mathbf{L}_k is a $M \times m_k$ full rank column matrix. For example if we observed the subset (y_{1k}, y_{3k}) only, then the $2 \times M$ matrix \mathbf{L}_k^T is given by

$$\mathbf{L}_k^T = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \end{bmatrix}.$$

In some situations we observe the sum of some components of \mathbf{y}_k . For example, the sum of all components is obtained by using the $1 \times M$ vector \mathbf{L}_k^T given by $\mathbf{L}_k^T = [1 \ \dots \ 1] \equiv \mathbf{I}_M^T$, where \mathbf{I}_M is the $M \times 1$ vector of 1's.

Note that $\mathbf{I}_k^{(o)}$ varies from one unit to another. The intention of the imputation process is to provide either a complete vector $\hat{\mathbf{y}}_k^{(e)} = (\hat{y}_{1k}^{(e)}, \dots, \hat{y}_{Mk}^{(e)})^T$ through conditional mean imputation or a complete vector $\hat{\mathbf{y}}_k^{(R)} = (\hat{y}_{1k}^{(R)}, \dots, \hat{y}_{Mk}^{(R)})^T$ through random imputation as estimate for \mathbf{y}_k given the observed information $\mathbf{I}_k^{(o)}$, with $\hat{y}_{ik}^{(e)} = \hat{y}_{ik}^{(R)} = y_{ik}$ if item i for unit k is observed explicitly or implicitly. Our imputed value of \mathbf{y}_k given $\mathbf{I}_k^{(o)}$ will of course be some function of $\mathbf{I}_k^{(o)}$, say $\mathbf{y}(\mathbf{I}_k^{(o)})$.

5.3. Conditional Mean Imputation for Item Nonresponse

Suppose that $(\mathbf{y}_k, \mathbf{I}_k^{(o)})$ are jointly distributed variables and we wish to predict the $M \times 1$ vector \mathbf{y}_k from the $m_k \times 1$ observed vector $\mathbf{I}_k^{(o)}$. One can in fact find the best estimator $\mathbf{y}^*(\mathbf{I}_k^{(o)})$ in the sense of minimizing MSE over all estimators. Except for the case of multivariate normal distributions, the conditional expectation $\mathbf{y}^*(\mathbf{I}_k^{(o)})$ could be a complicated nonlinear function of $\mathbf{I}_k^{(o)}$. So we restrict the class of estimators to the so-called linear estimators,

$$\mathbf{y}_L^{(e)}(\mathbf{I}_k^{(o)}) = \boldsymbol{\mu}_k + \mathbf{Q}_{\mathbf{y}_k \mathbf{I}_k^{(o)}} \mathbf{Q}_{\mathbf{I}_k^{(o)} \mathbf{I}_k^{(o)}}^{-1} (\mathbf{I}_k^{(o)} - \boldsymbol{\mu}_{\mathbf{I}_k^{(o)},k}) \equiv E_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}), \quad (5.1)$$

with

$$MSE[\mathbf{y}_L^{(e)}(\mathbf{I}_k^{(o)})] = Var_m(\mathbf{y}_k) - \mathbf{Q}_{\mathbf{y}_k \mathbf{I}_k^{(o)}} \mathbf{Q}_{\mathbf{I}_k^{(o)} \mathbf{I}_k^{(o)}}^{-1} \mathbf{Q}_{\mathbf{I}_k^{(o)} \mathbf{y}_k} \equiv Var_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}),$$

where $\boldsymbol{\mu}_k = E_m(\mathbf{y}_k)$, $\boldsymbol{\mu}_{\mathbf{I}_k^{(o)},k} = E_m(\mathbf{I}_k^{(o)})$, $\mathbf{Q}_{\mathbf{y}_k \mathbf{I}_k^{(o)}} = Cov_m(\mathbf{y}_k, \mathbf{I}_k^{(o)})$ and $\mathbf{Q}_{\mathbf{I}_k^{(o)} \mathbf{I}_k^{(o)}} = Cov_m(\mathbf{I}_k^{(o)}, \mathbf{I}_k^{(o)})$. As shown by Goldberger (1962), the linear estimator given by (5.1) is the best linear unbiased predictor of \mathbf{y}_k under the general linear model. The relative error, \mathbf{re}_k , of \mathbf{y}_k based on the knowledge $\mathbf{I}_k^{(o)}$ can be defined by the diagonal elements of the matrix

$$\mathbf{E}_k = [diag Var_m(\mathbf{y}_k)]^{-1} diag Var_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}).$$

If $\mathbf{I}_k^{(o)} = \emptyset$ then $E_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}) = E_m(\mathbf{y}_k)$, $Var_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}) = Var_m(\mathbf{y}_k)$ and $\mathbf{E}_k = \mathbf{I}_M$, where \mathbf{I}_M is the $M \times M$ identity matrix. On the other hand, if $\mathbf{I}_k^{(o)} = \mathbf{y}_k$ then $E_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}) = \mathbf{y}_k$, $Var_{mL}(\mathbf{y}_k | \mathbf{I}_k^{(o)}) = \mathbf{0}$ and $\mathbf{E}_k = \mathbf{0}_M$, where $\mathbf{0}_M$ is the $M \times M$ matrix of zero's.

Consider the case of $M = 2$, where $\mathbf{y}_k = (y_{1k}, y_{2k})^T$, $k = 1, \dots, N$, are independent observations generated from a bivariate normal distribution with means $\boldsymbol{\mu} = (\mu_1, \mu_2)^T$, variance $\boldsymbol{\sigma}^2 = (\sigma_1^2, \sigma_2^2)^T$ and coefficient of correlation $\rho_{y_1 y_2}$. We have $\sigma_t^2 = \sigma_1^2 + \sigma_2^2 + 2\rho_{y_1 y_2} \sigma_1 \sigma_2$ for $t = y_1 + y_2$, $Cov_m(y_i, t) = \sigma_i^2 + \rho_{y_1 y_2} \sigma_1 \sigma_2$ and $\rho_{y_i t}^2 = [\sigma_i + \rho_{y_1 y_2} \sigma_j]^2 / [\sigma_1^2 + \sigma_2^2 + 2\rho_{y_1 y_2} \sigma_1 \sigma_2]$ for $i \neq j$. Let $\alpha_{ij} = Cov_m(y_i, y_j) / Var_m(y_j) = \rho_{y_1 y_2} \sigma_i / \sigma_j$ for $i \neq j$, and $\alpha_{it} = Cov_m(y_i, t) / \sigma_t^2$. Table 2 reports the imputed values and the relative errors given some patterns of missingness.

Table 2: Imputed value, $\mathbf{y}_k^{(e)}$, and relative error given some response patterns when $\mathbf{y}_k \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Pattern of Missing Data			$E_m(\mathbf{y}_k \mathbf{I}_k^{(o)}, \boldsymbol{\beta})$		$\mathbf{re}(\mathbf{y}_k)$	
y_1	y_2	t	$y_1^{(e)}$	$y_2^{(e)}$	$re(y_1)$	$re(y_2)$
1	1	0	y_1	y_2	0	0
1	0	1	y_1	y_2	0	0
1	0	0	y_1	$\mu_2 + \alpha_{21}(y_{1k} - \mu_1)$	0	$1 - \rho_{y_1 y_2}^2$
0	0	1	$\mu_1 + \alpha_{1t}(t_k - \mu_t)$	$\mu_2 + \alpha_{2t}(t_k - \mu_t)$	$1 - \rho_{y_1 t}^2$	$1 - \rho_{y_2 t}^2$
0	0	0	μ_1	μ_2	1	1

The left panel of Table 3 shows all the true values related to Table 1. The middle panel shows the imputed values using conditional mean imputation, and the right panel shows the relative error in percentage due to imputation. The middle and right panels are obtained using our method studied in Demnati and Rao (2014), assuming the model parameter is known. Although the responses are generated independently with the same probability of response, it is seen from the right panel of Table 3 that correlations between variables makes approximate rather than fixed and exact response indicators compared to binary variable where the indicators take only zero or one. The unobserved item have a value that ranges in error between zero and one, and this value depends on the observed components and their relationship to the missing item. Although the probability of response is maintained the same for all variables, variable 3 has the lowest mean relative error.

Table 3 : True-, Imputed-values and Relative errors

Unit	t	True values			Imputed values			Relative errors (%)		
		y_1	y_2	y_3	$y_1^{(e)}$	$y_2^{(e)}$	$y_3^{(e)}$	$100re_1$	$100re_2$	$100re_3$
1	Y	7.12	19.71	34.13	7.12	19.71	34.13	0	0	0
2	Y	2.26	-13.53	-7.74	2.26	-13.53	-7.74	0	0	0
3	Y	5.83	12.93	20.01	5.83	13.21	19.73	0	5	2
4	Y	5.98	6.59	6.93	3.61	6.48	9.41	46	13	2
5	N	10.04	14.03	23.22	10.04	15.04	25.58	0	75	51
6	N	-0.8	-3.97	-7.23	5	10	15	100	100	100
	⋮									
Mean		3.54	10.87	12.89	4.14	11.36	15.80	20.6	15.5	13.35

Demnati and Rao (2014) developed estimators and variance estimators in the case of conditional multivariate imputation for item nonresponse.

REFERENCES

- Binder, D. A. (1983). "On the variances of asymptotically normal estimators from complex surveys". *International Statistical Review*, **51**, 279-292.
- Binder, D. A. (1996). "Linearization methods for single phase and two-phase samples: a cookbook approach". *Survey Methodology*, **22**, 17-22.
- Demnati, A. and Rao, J. N. K. (2002). "Linearization variance estimators for survey data with missing responses". *Proceeding of the Section Survey Research Methods*, American Statistical Association, 736-740.
- Demnati, A. and Rao, J.N.K. (2004). "Linearization variance estimators for survey data (with discussion)". *Survey Methodology*, **30**, 17-34.
- Demnati, A. and Rao, J.N.K. (2010). "Linearization variance estimators for model parameters from complex survey data". *Survey Methodology*, **36**, 193-199.
- Demnati, A. and Rao, J.N.K. (2014). "Multivariate conditional mean imputation for general pattern of missingness under control totals". In progress.
- Goldberger, A.S. (1962). "Best Linear unbiased prediction in the generalized linear regression model". *Journal of the American Statistical Association*, **57**, 369-375.

Royall, R. M. and Cumberland, W. G. (1981). "An empirical study of the ratio estimator and estimators of its variance". *Journal of the American Statistical Association*, **76**, 66-77.

Särndal, C.-E., Swensson, B. and Wretman, J.H. (1989). "The weighted residual technique for estimating the variance of the general regression estimator of the finite population total". *Biometrika*, **76**, 527-537.