

Combining information from multiple surveys through the empirical likelihood method

Changbao WU

Key words and phrases: Benchmark constraint; consistency requirement; generalized regression estimator; Newton–Raphson algorithm; pseudo empirical likelihood.

MSC 2000: Primary 62D05; secondary 62G09.

Abstract: It is often desirable to combine information collected in compatible multiple surveys in order to improve estimation and meet consistency requirements. Zieschang (1990) and Renssen & Nieuwenbroek (1997) suggested to this end the use of the generalized regression estimator with enlarged number of auxiliary variables. Unfortunately, adjusted weights associated with their approach can be negative. The author uses the notion of pseudo empirical likelihood to construct new estimators that are consistent, efficient and possess other attractive properties. The proposed approach is asymptotically equivalent to the earlier one, but it has clear maximum likelihood interpretations and its adjusted weights are always positive. The author also provides efficient algorithms for computing his estimators.

Regrouper l'information d'enquêtes multiples par la méthode de la vraisemblance empirique

Résumé : Il est souvent souhaitable de regrouper l'information de diverses enquêtes compatibles de façon à améliorer l'estimation et à assurer une certaine cohérence. Zieschang (1990) et Renssen & Nieuwenbroek (1997) ont suggéré à cette fin l'emploi d'un estimateur de régression généralisé exploitant un nombre accru de variables auxiliaires. Hélas, les poids ajustés liés à leur approche peuvent être négatifs. L'auteur tire de l'approche par la vraisemblance pseudo empirique de nouveaux estimateurs qui sont à la fois cohérents, efficaces et possèdent d'autres bonnes propriétés. L'approche proposée est asymptotiquement équivalente à la précédente mais a une interprétation claire en termes de vraisemblance maximale et ses poids ajustés sont toujours positifs. L'auteur fournit aussi des algorithmes efficaces pour le calcul de ses estimateurs.

1. INTRODUCTION

In survey practice, weight adjustment is routinely performed to accommodate, among other things, internal consistency requirements that are of interest both to survey statisticians and to the potential users of the survey data. Benchmark constraints are most commonly imposed where the adjusted weights w_i reproduce the known population totals (or means) of auxiliary variables \mathbf{x} , that is, $\sum_{i \in s} w_i \mathbf{x}_i = \mathbf{X}$, where s represents the set of sampled units and \mathbf{X} is the vector of known population totals. Such an adjustment can be achieved by using the generalized regression estimator. The generalized regression estimator is not only a vehicle to achieve the benchmark constraints; it is also more efficient when compared to the baseline Horvitz–Thompson estimator.

When two (or more) surveys are conducted for the same target population, another consistency requirement may arise. If some auxiliary variables are jointly collected in both surveys but their population totals are unknown, then it is desirable that, in addition to benchmark constraints over auxiliary variables with known population totals, the weights of both surveys produce the same estimates for the unknown population totals of the common auxiliary variables. This problem has previously been addressed by Zieschang (1990) and Renssen & Nieuwenbroek (1997). Both works proposed that the generalized regression estimator with an enlarged number of auxiliary variables be used to achieve that goal.

The generalized regression approach, however, has an undesirable property that was already being recognized by the authors. To quote Renssen & Nieuwenbroek (1997): “A disadvantage of the method is the increased possibility of negative weights, due to the enlarged number of

explanatory variables. The occurrence of negative weights is inherent to the general regression estimator, and for many users this is an undesirable feature.”

We propose to use the recently developed pseudo empirical likelihood (EL) method to construct estimators that not only meet the efficiency and consistency requirements but also have other attractive features. The EL method is a powerful nonparametric inference tool with applications in many areas of statistics. See Owen (2001) for a comprehensive account (and updated overview) of the subject. Historically, however, this method was first used in survey sampling by Hartley & Rao (1968). Its discrete and nonparametric nature is particularly appealing for finite population problems. In this article, we demonstrate that the EL approach is well suited to the current context, and consistency and efficiency requirements between two or multiple surveys can naturally be formed as constraints and can be integrated into the maximum likelihood estimation process. The two approaches, generalized regression and empirical likelihood (EL), are asymptotically equivalent but the latter has clear maximum likelihood interpretations and the resulting weights are always positive.

We consider two surveys in what follows, but our method can be extended to handle multiple surveys. A logically sound approach involves a joint maximum likelihood estimation using two samples. This is presented in Section 2. Also in Section 2, we present two algorithms for computing the proposed EL estimator. The first algorithm involves the profile likelihood method in searching for a solution and is efficient only when the common auxiliary variable is of dimension one. The second algorithm employs a novel reformulation of the problem and can easily be applied in general situations using the well developed algorithm of Chen, Sitter & Wu (2002). In Section 3, a separate empirical likelihood approach is employed where the EL estimators are computed separately for each survey with the unknown population means of the common auxiliary variables estimated from the combined sample data and used as control values. Computation in this case is simple and straightforward. The finite sample performance of the proposed EL estimators, with comparison to the generalized regression estimators of Zieschang (1990) and of Renssen & Nieuwenbroek (1997), is investigated in Section 4 through a simulation study. Variance estimation for the proposed estimators is discussed in Section 5. We conclude with a brief discussion on extending the method to multiple surveys and with some remarks in Section 6.

2. THE COMBINED EMPIRICAL LIKELIHOOD APPROACH

Suppose the finite population consists of N identifiable units. Associated with the i th unit are values of the study variables y_1 and y_2 and the vectors of auxiliary variables \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{z} , denoted by y_{1i} , y_{2i} , \mathbf{x}_{1i} , \mathbf{x}_{2i} , and \mathbf{z}_i , respectively, for $i = 1, \dots, N$. Information on (y_1, \mathbf{x}_1) is collected in the first survey and information on (y_2, \mathbf{x}_2) is gathered in the second survey. In addition, data on the common auxiliary variables \mathbf{z} are collected in both surveys. The two surveys, however, are carried out independently.

The two sets of sample data are $\{(y_{1i}, \mathbf{x}_{1i}, \mathbf{z}_i), i \in s_1\}$ and $\{(y_{2j}, \mathbf{x}_{2j}, \mathbf{z}_j), j \in s_2\}$, where s_1 and s_2 are the sets of sampled units from the first and the second survey, respectively. The population means $\bar{\mathbf{X}}_t = N^{-1} \sum_{i=1}^N \mathbf{x}_{ti}$ are known ($t = 1, 2$), but $\bar{\mathbf{Z}} = N^{-1} \sum_{i=1}^N \mathbf{z}_i$ are unknown. Zieschang (1990) and Renssen & Nieuwenbroek (1997) provided excellent motivations and real examples on this setting, including a highly valuable application on the split questionnaire survey designs. Let $\bar{Y}_t = N^{-1} \sum_{i=1}^N y_{ti}$, $t = 1, 2$, be the population quantities of interest. If y_1 and y_2 measure the same characteristic but over different time periods, then the difference $\Delta = \bar{Y}_2 - \bar{Y}_1$ may also be of interest.

Following arguments similar to those in Chen & Sitter (1999), the combined pseudo empirical log-likelihood function based on the two samples can be written as

$$\ell(\mathbf{p}, \mathbf{q}) = \sum_{i \in s_1} d_{1i} \log(p_i) + \sum_{j \in s_2} d_{2j} \log(q_j),$$

where $\mathbf{p} = (p_1, \dots, p_{n_1})'$, $\mathbf{q} = (q_1, \dots, q_{n_2})'$, $p_i = P(y_1 = y_{1i})$, $q_j = P(y_2 = y_{2j})$, $d_{ti} =$

$1/\pi_{ti}$, the π_{ti} are the first order inclusion probabilities, and n_t is the sample size for the t th survey, $t = 1, 2$.

The maximum pseudo empirical likelihood estimators for \bar{Y}_1 and \bar{Y}_2 are defined as

$$\hat{Y}_1 = \sum_{i \in s_1} \hat{p}_i y_{1i} \quad \text{and} \quad \hat{Y}_2 = \sum_{j \in s_2} \hat{q}_j y_{2j},$$

where the \hat{p}_i and \hat{q}_j , which are interpreted as the adjusted weights, maximize the joint pseudo empirical likelihood function $\ell(\mathbf{p}, \mathbf{q})$ subject to a system of normalization and consistency requirements:

$$\sum_{i \in s_1} p_i = 1 \quad (p_i > 0), \quad \sum_{j \in s_2} q_j = 1 \quad (q_j > 0), \quad (1)$$

$$\sum_{i \in s_1} p_i \mathbf{x}_{1i} = \bar{\mathbf{X}}_1, \quad \sum_{j \in s_2} q_j \mathbf{x}_{2j} = \bar{\mathbf{X}}_2, \quad (2)$$

$$\sum_{i \in s_1} p_i \mathbf{z}_i = \sum_{j \in s_2} q_j \mathbf{z}_j. \quad (3)$$

Both sets of benchmark constraints in (2) could involve measurements on the same \mathbf{x} variables and hence the same population means as well. In the absence of known population means, some or all of the equations in (2) can be removed from the system. The last set of equations (3) brings consistency between the two surveys over the common auxiliary variables. They also make the resulting estimators \hat{Y}_1 and \hat{Y}_2 more efficient by using the combined information from both surveys.

One of the related issues here is the existence of the foregoing defined combined EL estimators. The maximum pseudo empirical likelihood estimators \hat{Y}_t will not exist if $\bar{\mathbf{X}}_t$ is not an inner point of the convex hull formed by $\{\mathbf{x}_{ti}, i \in s_t\}$, or if the two convex hulls formed by $\{\mathbf{z}_i, i \in s_1\}$ and $\{\mathbf{z}_j, j \in s_2\}$ are disjoint. This occurs with probability approaching to zero as both sample sizes go to infinity. A proof of this can be sketched along the lines of Lemma 1 of Chen & Sitter (1999).

Another practically important issue is the computational aspect of the proposed EL method. We present two algorithms. Both of them take advantage of the well-behaved algorithm of Chen, Sitter & Wu (2002) for computing maximum empirical likelihood estimators under a single non-stratified sample. The first algorithm is efficient when the common auxiliary variable is univariate, while the second algorithm can be used in general situations.

2.1. The first algorithm.

Let $\sum_{i \in s_1} p_i \mathbf{z}_i = \sum_{j \in s_2} q_j \mathbf{z}_j = \boldsymbol{\theta}$ be fixed. It is then straightforward to show by the Lagrange multiplier method that

$$\hat{p}_i = \frac{d_{1i}^*}{1 + \boldsymbol{\lambda}'_1 \mathbf{u}_{1i}(\boldsymbol{\theta})}, \quad \hat{q}_j = \frac{d_{2j}^*}{1 + \boldsymbol{\lambda}'_2 \mathbf{u}_{2j}(\boldsymbol{\theta})}, \quad (4)$$

where $d_{ti}^* = d_{ti} / \sum_{i \in s_t} d_{ti}$ and

$$\mathbf{u}_{ti}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{x}_{ti} - \bar{\mathbf{X}}_t \\ \mathbf{z}_{ti} - \boldsymbol{\theta} \end{pmatrix},$$

with the understanding that \mathbf{z}_{ti} refers to \mathbf{z}_i from the t th sample, $t = 1, 2$. The Lagrange multipliers $\boldsymbol{\lambda}_1$ and $\boldsymbol{\lambda}_2$ are the solutions to

$$\sum_{i \in s_1} \frac{d_{1i}^* \mathbf{u}_{1i}(\boldsymbol{\theta})}{1 + \boldsymbol{\lambda}'_1 \mathbf{u}_{1i}(\boldsymbol{\theta})} = 0 \quad \text{and} \quad \sum_{j \in s_2} \frac{d_{2j}^* \mathbf{u}_{2j}(\boldsymbol{\theta})}{1 + \boldsymbol{\lambda}'_2 \mathbf{u}_{2j}(\boldsymbol{\theta})} = 0, \quad (5)$$

respectively. Then we obtain the profile likelihood function for θ by putting \hat{p}_i and \hat{q}_j into $\ell(\mathbf{p}, \mathbf{q})$ as given by (with a constant term omitted)

$$\ell(\theta) = - \sum_{i \in s_1} d_{1i} \log\{1 + \lambda'_1 \mathbf{u}_{1i}(\theta)\} - \sum_{j \in s_2} d_{2j} \log\{1 + \lambda'_2 \mathbf{u}_{2j}(\theta)\}.$$

The maximum point of $\ell(\theta)$, denoted by $\hat{\theta}$, can be found through the conventional profile analysis. We obtain the final adjusted weights \hat{p}_i and \hat{q}_j by plugging $\hat{\theta}$ and the associated λ_t into (4).

This algorithm involves finding λ_t ($t = 1, 2$) as solutions to (5) for each fixed value of θ , and then finding $\hat{\theta}$ that maximizes $\ell(\theta)$. For the first part, a simple and stable algorithm for solving (5) to obtain the vector-valued λ_t has been developed by Chen, Sitter & Wu (2002). As for $\hat{\theta}$, if the common auxiliary variable z is of dimension one, it can easily be found through the usual profile likelihood method. When z is high dimensional, so is θ , and this algorithm becomes awkward and impracticable. A more flexible algorithm is needed.

2.2. The second algorithm.

Suppose $\mathbf{z}_i = (z_{1i}, \dots, z_{ki})'$ is of dimension k . If we augment \mathbf{z}_i to have dimension $k + 1$ by including $z_{(k+1)i} = 1$ as the last component, we can rewrite the system of constraints (1)–(3) as

$$\sum_{i \in s_1} p_i + \sum_{j \in s_2} q_j = 2, \quad (6)$$

$$\begin{pmatrix} \mathbf{X}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^{(2)} \\ \mathbf{Z}^{(1)} & -\mathbf{Z}^{(2)} \end{pmatrix} \begin{pmatrix} \mathbf{p} \\ \mathbf{q} \end{pmatrix} = \begin{pmatrix} \overline{\mathbf{X}}_1 \\ \overline{\mathbf{X}}_2 \\ \mathbf{0} \end{pmatrix}, \quad (7)$$

where $\mathbf{X}^{(t)} = (\mathbf{x}_{t1}, \dots, \mathbf{x}_{tn_t})$, $\mathbf{Z}^{(t)} = (\mathbf{z}_{t1}, \dots, \mathbf{z}_{tn_t})$, \mathbf{z}_{ti} represents \mathbf{z}_i from the t th survey with 1 as its last component, $t = 1, 2$. Note that the very last equation in the system of (7) is $\sum_{i \in s_1} p_i - \sum_{j \in s_2} q_j = 0$; this together with (6) implies that $\sum_{i \in s_1} p_i = 1$ and $\sum_{j \in s_2} q_j = 1$.

We can further rewrite (7) as

$$\sum_{i \in s_1} p_i \mathbf{u}_{1i} + \sum_{j \in s_2} q_j \mathbf{u}_{2j} = \mathbf{0}, \quad (8)$$

where

$$\mathbf{u}_{1i} = \begin{pmatrix} \mathbf{x}_{1i} \\ \mathbf{0} \\ \mathbf{z}_{1i} \end{pmatrix} - \frac{1}{2} \begin{pmatrix} \overline{\mathbf{X}}_1 \\ \overline{\mathbf{X}}_2 \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{u}_{2j} = \begin{pmatrix} \mathbf{0} \\ \mathbf{x}_{2j} \\ -\mathbf{z}_{2j} \end{pmatrix} - \frac{1}{2} \begin{pmatrix} \overline{\mathbf{X}}_1 \\ \overline{\mathbf{X}}_2 \\ \mathbf{0} \end{pmatrix}. \quad (9)$$

It is now clear that maximizing $\ell(\mathbf{p}, \mathbf{q})$ under the restrictions (1), (2) and (3) is equivalent to maximizing $\ell(\mathbf{p}, \mathbf{q})$ subject to (6) and (8). By using the Lagrange multiplier method, we can show that

$$\hat{p}_i = \frac{d_{1i}^*}{1 + \lambda'_1 \mathbf{u}_{1i}}, \quad \hat{q}_j = \frac{d_{2j}^*}{1 + \lambda'_2 \mathbf{u}_{2j}},$$

where

$$d_{ti}^* = 2d_{ti} / \left(\sum_{i \in s_1} d_{1i} + \sum_{j \in s_2} d_{2j} \right)$$

for $t = 1, 2$, and the common Lagrange multiplier λ is the solution to

$$\sum_{t=1,2} \sum_{i \in s_t} \frac{d_{ti}^* \mathbf{u}_{ti}}{1 + \lambda'_t \mathbf{u}_{ti}} = \mathbf{0}. \quad (10)$$

The modified Newton–Raphson algorithm of Chen, Sitter & Wu (2002) can ideally be used to solve (10). Although such a modification is necessary for the theoretical proof of convergence, it is our experience that the following conventional Newton–Raphson iteration procedure works well for almost all cases:

$$\boldsymbol{\lambda}^{(m+1)} = \boldsymbol{\lambda}^{(m)} + \left\{ \sum_{t=1,2} \sum_{i \in s_t} \frac{d_{ti}^* \mathbf{u}_{ti} \mathbf{u}_{ti}'}{(1 + [\boldsymbol{\lambda}^{(m)}]' \mathbf{u}_{ti})^2} \right\}^{-1} \sum_{t=1,2} \sum_{i \in s_t} \frac{d_{ti}^* \mathbf{u}_{ti}}{1 + [\boldsymbol{\lambda}^{(m)}]' \mathbf{u}_{ti}},$$

with the initial value of $\boldsymbol{\lambda}$ chosen as $\mathbf{0}$.

This second algorithm is applicable in general situations. It requires solving (10) only once using the existing well-behaved algorithm of Chen, Sitter & Wu (2002) and can be programmed by survey users with popular statistical software such as SAS or R/S-PLUS.

2.3. A comparison with Zieschang's regression method.

The combined empirical likelihood approach proposed in this article is in the same spirit as the composite generalized regression estimator proposed by Zieschang (1990). This is evident when we compare the constraints (7) used here with the enlarged regression system (3.10) used by Zieschang. There are several advantages, however, in using the empirical likelihood method. In addition to its clear maximum likelihood interpretations, we compute the EL estimator based on the normalized intrinsically positive weights, that is, $\hat{p}_i > 0$ and $\sum_{i \in s_1} \hat{p}_i = 1$. This latter feature is particularly appealing to potential users of the survey data since the published weights are often used for a variety of purposes, including the estimation of proportions or more generally the finite population distribution function $F(y)$. The EL estimator $\hat{F}_{EL}(y)$ itself will be a genuine distribution function. It is range-respecting and can be inverted directly to get quantile estimates.

An explicit relationship between the EL estimator and a generalized regression-type estimator can be established.

THEOREM 1. *Under suitable regularity conditions, the maximum pseudo empirical likelihood estimators*

$$\hat{Y}_1 = \sum_{i \in s_1} \hat{p}_i y_{1i} \quad \text{and} \quad \hat{Y}_2 = \sum_{j \in s_2} \hat{q}_j y_{2j}$$

are asymptotically equivalent to a generalized regression-type estimator, that is,

$$\hat{Y}_t = \bar{y}_t + \hat{\mathbf{B}}'_{t1} (\bar{\mathbf{X}}_1 - \bar{\mathbf{x}}_1) + \hat{\mathbf{B}}'_{t2} (\bar{\mathbf{X}}_2 - \bar{\mathbf{x}}_2) + \hat{\mathbf{B}}'_{t3} (\bar{\mathbf{z}}_2 - \bar{\mathbf{z}}_1) + o_p(n^{-1/2}), \quad (11)$$

where

$$\bar{y}_t = \sum_{i \in s_t} d_{ti}^* y_{ti}, \quad \bar{\mathbf{x}}_t = \sum_{i \in s_t} d_{ti}^* \mathbf{x}_{ti}, \quad \bar{\mathbf{z}}_t = \sum_{i \in s_t} d_{ti}^* \mathbf{z}_{ti}, \quad n = n_1 + n_2$$

and the combined "regression coefficients" $\hat{\mathbf{B}}_t = (\hat{\mathbf{B}}'_{t1}, \hat{\mathbf{B}}'_{t2}, \hat{\mathbf{B}}'_{t3})'$ are computed as

$$\hat{\mathbf{B}}_t = \left(\sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} \mathbf{u}_{ti}' \right)^{-1} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} y_{ti},$$

with the \mathbf{u}_{ti} defined by (9).

The required regularity conditions and a proof of the theorem are given in the Appendix. Note that the combined auxiliary information $\bar{\mathbf{X}}_1$, $\bar{\mathbf{X}}_2$, $\bar{\mathbf{z}}_1$ and $\bar{\mathbf{z}}_2$, as well as the basic design weights d_{1i} and d_{2j} from both surveys all appear explicitly in the equivalent generalized regression estimator, an estimator that is quite unique from the conventional point of view. Further, if both sampling designs satisfy $\sum_{i \in s_1} d_{1i} = \sum_{i \in s_2} d_{2i} = N$, as is the case under simple random sampling or stratified random sampling, then $d_{ti}^* = d_{ti}/N$, and the estimators \bar{y}_t , $\bar{\mathbf{x}}_t$ and $\bar{\mathbf{z}}_t$ all reduce to the usual Horvitz–Thompson estimators for the corresponding population means.

3. THE SEPARATE EMPIRICAL LIKELIHOOD APPROACH

In the combined approach, the unknown population mean vector $\bar{\mathbf{Z}}$ is implicitly estimated by the maximum pseudo empirical likelihood estimator $\hat{\boldsymbol{\theta}}$ from the pooled sample. This can be seen from the first algorithm presented in Section 2.1. Some computational complications arising from the combined approach are due solely to the attempt to estimate $\bar{\mathbf{Z}}$ by $\hat{\boldsymbol{\theta}}$.

One way to circumvent this difficulty is to take a two-step approach. Suppose we replace $\hat{\boldsymbol{\theta}}$ by a different estimator of $\bar{\mathbf{Z}}$, say $\bar{\mathbf{z}}$, using the combined data from both surveys. We then use the entries in this $\bar{\mathbf{z}}$ as control values for the constraints used in the empirical likelihood estimation for each of the two surveys. By doing so, we not only bring consistency for the common auxiliary variables \mathbf{z} between the two surveys but also improve the resulting estimators \hat{Y}_1 and \hat{Y}_2 if $\bar{\mathbf{z}}$ is suitably constructed from the combined sample data. This is similar to the case of two-phase sampling where the unknown population quantity $\bar{\mathbf{Z}}$ is estimated using the large first-phase sample.

The estimation of \bar{Y}_1 and \bar{Y}_2 with a pre-determined $\bar{\mathbf{z}}$ as control value becomes two separate estimation problems. For instance, the EL estimator for \bar{Y}_1 is given by $\hat{Y}_1 = \sum_{i \in s_1} \hat{p}_i y_{1i}$, where the \hat{p}_i maximize $\ell(\mathbf{p}) = \sum_{i \in s_1} d_{1i} \log(p_i)$ subject to constraints

$$\sum_{i \in s_1} p_i = 1 \quad (p_i > 0), \quad \sum_{i \in s_1} p_i \mathbf{x}_{1i} = \bar{\mathbf{X}}_1 \quad \text{and} \quad \sum_{i \in s_1} p_i \mathbf{z}_{1i} = \bar{\mathbf{z}}.$$

The resulting weights are computed as $\hat{p}_i = d_{1i}^* / (1 + \boldsymbol{\lambda}' \mathbf{u}_{1i})$, where $d_{1i}^* = d_{1i} / \sum_{i \in s_1} d_{1i}$ and the Lagrange multiplier $\boldsymbol{\lambda}$ is the solution to

$$\sum_{i \in s_1} \frac{d_{1i}^* \mathbf{u}_{1i}}{1 + \boldsymbol{\lambda}' \mathbf{u}_{1i}} = 0, \quad \text{with} \quad \mathbf{u}_{1i} = \begin{pmatrix} \mathbf{x}_{1i} - \bar{\mathbf{X}}_1 \\ \mathbf{z}_{1i} - \bar{\mathbf{z}} \end{pmatrix}. \quad (12)$$

The algorithm of Chen, Sitter & Wu (2002) can be used directly here to obtain $\boldsymbol{\lambda}$ without any modification.

The major issue in this separate EL approach is the choice of $\bar{\mathbf{z}}$. Renssen & Nieuwenbroek (1997) provided an excellent account of the estimation of $\bar{\mathbf{Z}}$ using combined sample data. They suggested a general class of estimators of the form $\bar{\mathbf{z}} = \mathbf{P}\bar{\mathbf{z}}_1 + \mathbf{Q}\bar{\mathbf{z}}_2$, where \mathbf{P} and \mathbf{Q} are two matrices with compatible dimensions such that $\mathbf{P} + \mathbf{Q} = \mathbf{I}$, and $\bar{\mathbf{z}}_t$ is the generalized regression estimator of \mathbf{Z} with \mathbf{x}_t as auxiliary variables. In the absence of $\bar{\mathbf{X}}_t$, one can take $\bar{\mathbf{z}}_t$ as the Horvitz–Thompson estimator of $\bar{\mathbf{Z}}$ using data from the t th survey.

Two choices of the matrix pair (\mathbf{P}, \mathbf{Q}) will be examined in the simulation study presented in the next section. The simplest one is the proportional combination, where $\mathbf{P} = (n_1 + n_2)^{-1} n_1 \mathbf{I}$ and $\mathbf{Q} = (n_1 + n_2)^{-1} n_2 \mathbf{I}$; the optimal combination uses

$$\mathbf{P} = V_p(\bar{\mathbf{z}}_2) \{V_p(\bar{\mathbf{z}}_1) + V_p(\bar{\mathbf{z}}_2)\}^{-1} \quad \text{and} \quad \mathbf{Q} = V_p(\bar{\mathbf{z}}_1) \{V_p(\bar{\mathbf{z}}_1) + V_p(\bar{\mathbf{z}}_2)\}^{-1},$$

where $V_p(\bar{\mathbf{z}}_t)$ is the design-based variance-covariance matrix of $\bar{\mathbf{z}}_t$. Note that the \mathbf{P} used in the optimal combination can also be written as

$$\mathbf{P} = [\{V_p(\bar{\mathbf{z}}_1)\}^{-1} + \{V_p(\bar{\mathbf{z}}_2)\}^{-1}]^{-1} \{V_p(\bar{\mathbf{z}}_1)\}^{-1},$$

and similarly for \mathbf{Q} as well. This choice is optimal since it minimizes $V_p(\mathbf{a}'\bar{\mathbf{z}})$ for an arbitrary constant vector \mathbf{a} among the general class of estimators considered by Renssen & Nieuwenbroek (1997). When simple random sampling is used for both surveys and the $\bar{\mathbf{z}}_t$ are the simple sample means, the optimal combination reduces to the proportional one if the two sampling fractions are the same or can be ignored. Note that for the optimal combination, the matrices \mathbf{P} and \mathbf{Q} need to be replaced by sample-based estimates for applications.

The separate EL approach is less elegant than the combined one in terms of maximum likelihood estimation. This approach, however, is intuitively attractive, and computation in this case is straightforward and simple. Under suitable regularity conditions similar to those used in Theorem 1, we can show that the separate EL estimator is asymptotically equivalent to the regression estimator discussed by Renssen & Nieuwenbroek (1997), that is,

$$\widehat{Y}_t = \bar{y}_t + \widehat{B}'_{t1}(X_t - \bar{x}_t) + \widehat{B}'_{t2}(\bar{z} - \bar{z}_t^*) + o_p(n^{-1/2}), \quad (13)$$

where

$$\bar{y}_t = \sum_{i \in s_t} d_{ti}^* y_{ti}, \quad \bar{x}_t = \sum_{i \in s_t} d_{ti}^* x_{ti}, \quad \bar{z}_t^* = \sum_{i \in s_t} d_{ti}^* z_{ti}, \quad d_{ti}^* = d_{ti} / \sum_{i \in s_t} d_{ti},$$

and the regression coefficients $\widehat{B}_t = (\widehat{B}'_{t1}, \widehat{B}'_{t2})'$ are given by

$$\widehat{B}_t = \left(\sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} \mathbf{u}'_{ti} \right)^{-1} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} y_{ti},$$

where \mathbf{u}_{1i} (and \mathbf{u}_{2i} in obvious form) are defined in (12). The first two terms on the right-hand side of (13) can be viewed as a generalized regression estimator for \bar{Y}_t based on auxiliary variables \mathbf{x}_t and the third one is an adjusting term in an attempt to further improve the regression estimator with the extra information on \mathbf{z} variables.

It is worthwhile to note that for the separate EL estimator, the different sample sizes can easily be taken into account for the estimation of \bar{Z} . The combined EL approach, however, does not automatically accommodate this and requires a special weighting adjustment to achieve the same goal. Further development in this direction will not be pursued here, but this argument provides a possible explanation as to why the combined EL estimator is often outperformed by the separate one, as shown by the simulation results reported in the next section.

4. SIMULATION STUDY

In this section, we examine the finite sample performance of proposed estimators through a limited simulation study. The finite population used in this study was based on real data from Statistics Canada's 1996 Family Expenditure (FAMEX) Survey for the province of Ontario. The data set contains $N = 2396$ observations measured over a variety of characteristics. Variables which are relevant to our study include x_1 : number of children (age < 15); x_2 : number of youths (age 15–24); x_3 : number of people in the household; z : total income after taxes; and y : total expenditure.

In the simulation, we treat the data set itself as a finite population. This population is further split into eight strata according to the original sampling design. For the first survey, the number of children (x_1) and the number of people (x_3) with known population means are used as control variables, and the total expenditure y is treated as the response variable. We also assume that the variables x_2 and x_3 are used as control variables in the second survey, and information on total income (z) is conveniently collected for both surveys, but the population mean \bar{Z} is unknown. The goal is to estimate the population mean \bar{Y} using all useful information while respecting the consistency requirement imposed over the z variable for the two surveys.

For each simulation run, a stratified random sample of size n_t under proportional allocation is taken for the t th survey, $t = 1, 2$, and three maximum pseudo empirical likelihood estimators for \bar{Y} are computed. Let EL(C) denote the combined EL estimator, let EL(SP) be the separate EL estimator using proportional combination for the estimation of \bar{Z} , and let EL(SO) represent the separate EL estimator using optimal combination in estimating \bar{Z} . Also computed for each simulation are three generalized regression-type estimators: the one proposed by Zieschang (1990) is denoted by GR(Z), which is equivalent to our combined EL estimator EL(C); and the estimators

proposed by Renssen & Nieuwenbroek (1997) are denoted by GR(RN1) and GR(RN2), corresponding to our EL(SP) and EL(SO), respectively. The Λ matrix used to formulate the GR(Z) estimator is taken as $\text{diag}(d_1, \dots, d_n)$. The process is repeated independently $B = 1000$ times.

The performance of an estimator \widehat{Y} is measured in terms of the simulated Relative Bias (RB) and Relative Efficiency (RE) defined as

$$\text{RB} = \frac{1}{B} \sum_{b=1}^B \frac{\widehat{Y}(b) - \bar{Y}}{\bar{Y}} \quad \text{and} \quad \text{RE} = \frac{\text{MSE}(\widehat{Y}_0)}{\text{MSE}(\widehat{Y})},$$

where $\widehat{Y}(b)$ is the estimator \widehat{Y} computed from the b th simulated sample,

$$\text{MSE}(\widehat{Y}) = \frac{1}{B} \sum_{b=1}^B \{\widehat{Y}(b) - \bar{Y}\}^2,$$

and \widehat{Y}_0 is the baseline estimator for comparison. In our study, we have chosen \widehat{Y}_0 as the generalized regression estimator (GREG) of \bar{Y} using x_1 and x_3 as auxiliary variables. Note that the sample information on z cannot be used here for the regression estimation of \bar{Y} , since \bar{Z} is assumed to be unknown. The sample sizes $n_t = 80, 160$ and 240 used in the simulation represent a typical sampling fractions of 2.5%, 5% and 10% respectively.

TABLE 1: Simulated relative efficiencies based on the 1996 Statcan FAMEX survey data.

n_1	n_2	GREG	EL(C)	GR(Z)	EL(SP)	GR(RN1)	EL(SO)	GR(RN2)
80	80	1.00	1.19	1.13	1.28	1.28	1.28	1.28
	160	1.00	1.36	1.28	1.31	1.43	1.31	1.42
	240	1.00	1.42	1.31	1.43	1.49	1.48	1.49
160	80	1.00	1.03	0.91	1.15	1.15	1.15	1.15
	160	1.00	1.28	1.10	1.27	1.26	1.27	1.27
	240	1.00	1.33	1.18	1.29	1.28	1.30	1.30
240	80	1.00	0.91	0.80	1.16	1.14	1.15	1.13
	160	1.00	1.11	0.95	1.16	1.14	1.16	1.15
	240	1.00	1.24	1.07	1.22	1.20	1.24	1.23

The absolute values of the simulated relative biases are all less than 0.1% and are not reported here. Table 1 reports the relative efficiency of the EL and the generalized regression-type estimators under various scenarios for the sample size combinations. Our major findings can be summarized as follows:

- (i) The two separate EL estimators have similar performance and they both perform well. The gain of efficiency is more pronounced when the second sample size is larger.
- (ii) The combined EL estimator has satisfactory performance when the second sample has a comparable size but could have deteriorated performance otherwise (that is, the case of $n_1 = 240$ and $n_2 = 80$).
- (iii) The generalized regression estimators of Renssen & Nieuwenbroek (1997) have very similar performance to the separate EL estimators, but the generalized regression estimator of Zieschang (1990) is outperformed by the combined EL estimator in all cases.
- (iv) The use of information on the common auxiliary variable z provides substantial improvement over the baseline generalized regression estimator when the second sample is not too small.

5. VARIANCE ESTIMATION

The combined and the separate maximum pseudo empirical likelihood estimators developed in this article belong to the general class of nonlinear estimators and their exact design-based variances do not have a closed form. Simple consistent variance estimators, however, can be derived by using the asymptotically equivalent regression-type estimators given by (11) and (13). For the combined EL estimator, one complication arises due to the d_{ti}^* given by

$$2d_{ti} / \left(\sum_{i \in s_1} d_{1i} + \sum_{j \in s_2} d_{2j} \right).$$

Under the assumed regularity conditions, we have

$$N^{-1} \sum_{i \in s_t} d_{ti} = 1 + O_p(n^{-1/2})$$

for $t = 1, 2$, where $n = n_1 + n_2$. By using the delta method, we can approximate the nonlinear terms in (11) by first order linear expansions. For instance,

$$\bar{y}_1 = \sum_{i \in s_1} d_{1i}^* y_{1i} = \frac{1}{N} \sum_{i \in s_1} d_{1i} y_{1i} - \bar{Y}_1 \left(\frac{1}{N} \sum_{i \in s_1} d_{1i} - 1 \right) - \bar{Y}_1 \left(\frac{1}{N} \sum_{j \in s_2} d_{2j} - 1 \right) + o_p(n^{-1/2}).$$

It is straightforward to show that

$$\widehat{\bar{Y}}_1 = C + \frac{1}{N} \sum_{i \in s_1} d_{1i} A_{1i} + \frac{1}{N} \sum_{j \in s_2} d_{2j} A_{2j} + o_p(n^{-1/2}), \quad (14)$$

where C is a constant,

$$\begin{aligned} A_{1i} &= (y_{1i} - \mathbf{B}'_{11} \mathbf{x}_{1i} - \mathbf{B}'_{13} \mathbf{z}_{1i}) - (\bar{Y}_1 - \mathbf{B}'_{11} \bar{\mathbf{X}}_1 - \mathbf{B}'_{12} \bar{\mathbf{X}}_2), \\ A_{2j} &= (\mathbf{B}'_{13} \mathbf{z}_{2j} - \mathbf{B}'_{12} \mathbf{x}_{2j}) - (\bar{Y}_1 - \mathbf{B}'_{11} \bar{\mathbf{X}}_1 - \mathbf{B}'_{12} \bar{\mathbf{X}}_2), \end{aligned}$$

and \mathbf{B} is the combined population regression coefficients given by

$$\mathbf{B}_t = (\mathbf{B}'_{t1}, \mathbf{B}'_{t2}, \mathbf{B}'_{t3})' = \left(\sum_{t=1}^2 \sum_{i=1}^N \mathbf{u}_{ti} \mathbf{u}'_{ti} \right)^{-1} \sum_{i=1}^N \mathbf{u}_{ti} y_{ti}.$$

Since the two surveys are independent, the asymptotic variance of $\widehat{\bar{Y}}_1$ and a consistent variance estimator can be developed in obvious way by using (14) and the standard variance estimation method for the Horvitz–Thompson estimator of population means. Similar developments can also be made for the separate EL estimators. The details are omitted.

6. CONCLUDING REMARKS

The EL estimators developed in this article can be extended to handle multiple surveys. Suppose there are three sets of sample data $\{(y_{ti}, \mathbf{x}_{ti}, \mathbf{z}_{ti}), i \in s_t\}$, $t = 1, 2, 3$, where the population means $\bar{\mathbf{X}}_t = N^{-1} \sum_{i=1}^N \mathbf{x}_{ti}$ are known for $t = 1, 2, 3$ but the population means $\bar{\mathbf{Z}}$ for the common variables \mathbf{z} are unknown. For the separate EL approach, one can handle this by simply using $\bar{\mathbf{z}} = \mathbf{P}\bar{\mathbf{z}}_1 + \mathbf{Q}\bar{\mathbf{z}}_2 + \mathbf{R}\bar{\mathbf{z}}_3$ to estimate $\bar{\mathbf{Z}}$ under a suitable combination of \mathbf{P} , \mathbf{Q} and \mathbf{R} .

As for the combined approach, some modification is needed to deal with the computational complexities. The combined EL estimator for $\bar{Y}_t = N^{-1} \sum_{i=1}^N y_{ti}$ is given by $\widehat{\bar{Y}}_t = \sum_{i \in s_t} \hat{p}_{ti} y_{ti}$, where the \hat{p}_{ti} maximize

$$\ell(\mathbf{p}) = \sum_{t=1}^3 \sum_{i \in s_t} d_{ti} \log(p_{ti}),$$

subject to

$$\begin{aligned} \sum_{i \in s_t} p_{ti} &= 1 \quad (p_{ti} > 0), & \sum_{i \in s_t} p_{ti} \mathbf{x}_{ti} &= \bar{\mathbf{X}}_t, \quad t = 1, 2, 3, \\ \sum_{i \in s_1} p_{1i} \mathbf{z}_{1i} &= \sum_{i \in s_2} p_{2i} \mathbf{z}_{2i} = \sum_{i \in s_3} p_{3i} \mathbf{z}_{3i}. \end{aligned}$$

Modification of the first algorithm for the current situation is straightforward: we let

$$\sum_{i \in s_1} p_{1i} \mathbf{z}_{1i} = \sum_{i \in s_2} p_{2i} \mathbf{z}_{2i} = \sum_{i \in s_3} p_{3i} \mathbf{z}_{3i} = \boldsymbol{\theta}$$

be fixed first and then find $\boldsymbol{\theta}$ for the EL estimator through the profile likelihood method. The second algorithm can also be modified in this case. Let \mathbf{z}_{ti} be augmented to include 1 as its last component. Finding \hat{p}_{ti} is equivalent to maximizing $\ell(\mathbf{p})$ under the constraints

$$\sum_{t=1}^3 \sum_{i \in s_t} p_{ti} = 3, \quad \sum_{t=1}^3 \sum_{i \in s_t} p_{ti} \mathbf{u}_{ti} = \mathbf{0},$$

where

$$\mathbf{u}_{1i} = \begin{pmatrix} \mathbf{x}_{1i} \\ \mathbf{0} \\ \mathbf{0} \\ z_{1i} \\ z_{1i} \end{pmatrix} - \frac{1}{3} \begin{pmatrix} \bar{\mathbf{X}}_1 \\ \bar{\mathbf{X}}_2 \\ \bar{\mathbf{X}}_3 \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{u}_{2i} = \begin{pmatrix} \mathbf{0} \\ \mathbf{x}_{2i} \\ \mathbf{0} \\ -z_{2i} \\ \mathbf{0} \end{pmatrix} - \frac{1}{3} \begin{pmatrix} \bar{\mathbf{X}}_1 \\ \bar{\mathbf{X}}_2 \\ \bar{\mathbf{X}}_3 \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{u}_{3i} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{x}_{3i} \\ \mathbf{0} \\ -z_{3i} \end{pmatrix} - \frac{1}{3} \begin{pmatrix} \bar{\mathbf{X}}_1 \\ \bar{\mathbf{X}}_2 \\ \bar{\mathbf{X}}_3 \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}.$$

The final solution is given by $\hat{p}_{ti} = d_{ti}^*/(1 + \boldsymbol{\lambda}' \mathbf{u}_{ti})$, where $d_{ti}^* = 3d_{ti}/\sum_{t=1}^3 \sum_{i \in s_t} d_{ti}$ and $\boldsymbol{\lambda}$ is the solution to (10) but replacing the first summation with the one over all three samples.

Adjusting weights to satisfy certain efficiency and consistency requirements is a constant theme in survey sampling, and having positive adjusted weights is a highly desirable property for the users of production micro data files, where the weights are viewed as the number of units in the finite population represented by the sampled unit. Positive weights will also guarantee positive estimation for known positive population quantities.

It should be noted that positive weights in regression estimation can theoretically be achieved through constrained minimization under the context of calibration estimation as discussed in Deville & Särndal (1992). The practical implementation of such a method is not straightforward, however, and it often involves ad hoc approximations. The loss of efficiency due to these approximations is usually unknown. The empirical likelihood method, on the other hand, provides a natural way of doing this with final adjusted weights that are intrinsically positive. It should also be noted that the total amount of time required for computing the proposed EL estimators remains limited. In our simulation study, it takes less than 20 seconds on a dual process Sun Unix workstation to compute the combined EL estimator when $n_1 = n_2 = 240$ and the program is written in R/S-PLUS.

While the adjusted weights using a generalized regression-type technique tend to have some small or negative values, the weights obtained from the EL approach can occasionally contain a few large values. For the simulation results reported in Section 4, where a proportional sample size allocation scheme is used to draw the two stratified random samples, the g -weights given by $g_i = w_i/d_i$ are all within the range of (0.25, 4.00), where the w_i denote the EL adjusted weights and the d_i represent the basic design weights. More than 99% of these g -weights are indeed between 0.50 and 2.00. If we use an unbalanced allocation scheme, where the largest

stratum ($N_h = 763$) and the smallest one ($N_h = 33$) receive equal sample sizes, we observed that a couple of g -weights can be larger than 4.00 or even 6.00. Theoretically this is not a problem regarding the statistical properties of the EL estimators. For users who also have concerns about large weights, the idea of minimum relaxation of constraints presented in Chen, Sitter & Wu (2002) can be applied to obtain more general range restricted weights through the EL method.

It is easy to argue, both theoretically and empirically, that under ideal situations, the gain in efficiency from using the combined sample data is almost guaranteed. Cases where the forced consistency requirement over the common variables will likely be detrimental to the resulting estimators include: (1) severe uncontrolled nonsampling errors; (2) extremely unbalanced sampling designs or sample size allocations; (3) misconceptualized target populations; (4) weak or lack of correlation between the common variables and the response variables; and (5) use of questionable common variables.

The successful use of the proposed EL method for combining information on common auxiliary variables for real surveys requires detailed consideration at the planning stage and careful discretion at the estimation stage. As pointed out by Renssen & Nieuwenbroek (1997), common variables in the strict sense are not easily found due to discrepancies between definitions, methods of observation, and reference periods. Such complications can be reduced if the involved surveys are harmonized at the design stage. For example, in split questionnaire design where certain common questions are contained in both versions of the questionnaire, attention should be given to the ordering and positioning of these common questions to reduce the potential response bias and/or carry-over effect. For human population surveys conducted regularly over time, variables related to gender, age, educational background, etc., can easily be conceived as common variables when the change of population dynamics over a certain time period can be ignored. Other variables need to be treated with care: for instance "Employment Status in May 2003" can be measured as direct answers in a June 2003 survey, but for surveys conducted at a later time, only as recalled answers. Only if the recalled answers are as accurate as the direct ones can measurements of this type be treated as common variables.

APPENDIX

We assume that the maximum pseudo empirical likelihood estimators exist. A certain asymptotic framework is also needed at this moment. We refer the reader to Isaki & Fuller (1982) for details. The required regularity conditions are stated in terms of the basic design weights d_{ti} and the values of the auxiliary variables. The notation $\bar{\mathbf{x}}_t$, $\bar{\mathbf{z}}_t$ and d_{ti}^* follows from Section 2.2. Stochastic orders involving random vectors or matrices are interpreted as componentwise.

- (i) $\bar{\mathbf{X}}_t - \bar{\mathbf{x}}_t = O_p(n_t^{-1/2})$ and $\bar{\mathbf{Z}} - \bar{\mathbf{z}}_t = O_p(n_t^{-1/2})$;
- (ii) $\sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} \mathbf{u}'_{ti} - N^{-1} \sum_{t=1,2} \sum_{i=1}^N \mathbf{u}_{ti} \mathbf{u}'_{ti} = O_p(n^{-1/2})$;
- (iii) $|\sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} \mathbf{u}'_{ti}| \neq 0$ and $N^{-1} \sum_{t=1,2} \sum_{i=1}^N \mathbf{u}_{ti} \mathbf{u}'_{ti} = O(1)$.

Conditions (i) and (ii) are standard; the second part of condition (iii) states that the \mathbf{x} and the \mathbf{z} variables have finite second moments, and this implies

$$\max\{\mathbf{x}_{ti}, i \in s_t\} = o_p(n_t^{1/2}) \quad \text{and} \quad \max\{\mathbf{z}_{ti}, i \in s_t\} = o_p(n_t^{1/2})$$

(Owen 1990, Lemma 3). A key result for the proof of the theorem is to show $\boldsymbol{\lambda} = O_p(n^{-1/2})$. This can be (nonrigorously) argued as follows: by rewriting $d_{ti}^* \mathbf{u}_{ti}$ as $d_{ti}^* \{\mathbf{u}_{ti}(1 + \boldsymbol{\lambda}' \mathbf{u}_{ti}) - \mathbf{u}_{ti} \mathbf{u}'_{ti} \boldsymbol{\lambda}\}$ in (10), we have

$$\sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} = \left(\sum_{t=1,2} \sum_{i \in s_t} \hat{w}_{ti} \mathbf{u}_{ti} \mathbf{u}'_{ti} \right) \boldsymbol{\lambda},$$

where $\hat{w}_{1i} = \hat{p}_i$ and $\hat{w}_{2i} = \hat{q}_i$. Since

$$\sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} = [(\bar{\mathbf{x}}_1 - \bar{\mathbf{X}}_1)', (\bar{\mathbf{x}}_2 - \bar{\mathbf{X}}_2)', (\bar{\mathbf{z}}_1 - \bar{\mathbf{z}}_2)']'$$

is of order $O_p(n^{-1/2})$ (componentwise), $\sum_{t=1,2} \sum_{i \in s_t} \hat{w}_{ti} \mathbf{u}_{ti} \mathbf{u}_{ti}'$ is the maximum empirical likelihood estimator of $N^{-1} \sum_{t=1,2} \sum_{i=1}^N \mathbf{u}_{ti} \mathbf{u}_{ti}'$ which is of order $O(1)$, then we must have $\boldsymbol{\lambda} = O_p(n^{-1/2})$. A more rigorous argument leading to this conclusion will follow the lines of (2.11)–(2.14) of Owen (1990, pp. 100–101). It follows that $\max_{i \in s_t} |\boldsymbol{\lambda}' \mathbf{u}_{ti}| = o_p(1)$ and $(1 + \boldsymbol{\lambda}' \mathbf{u}_{ti})^{-1} = 1 - \boldsymbol{\lambda}' \mathbf{u}_{ti} \{1 + o_p(1)\}$, with the uniform term $o_p(1)$ over $i \in s_t$, $t = 1, 2$. Applying this last expansion to (10), we get

$$\boldsymbol{\lambda} = \left(\sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} \mathbf{u}_{ti}' \right)^{-1} \sum_{t=1,2} \sum_{i \in s_t} d_{ti}^* \mathbf{u}_{ti} + o_p(n^{-1/2}).$$

The final expansion for \hat{Y}_1 (or \hat{Y}_2) is a simple consequence of $\hat{p}_i = d_{1i}^* [1 - \boldsymbol{\lambda}' \mathbf{u}_{1i} \{1 + o_p(1)\}]$, where the term $o_p(1)$ is independent of i .

ACKNOWLEDGEMENTS

This research was supported by a grant from the Natural Sciences and Engineering Research Council of Canada. Comments and suggestions from Professor Jiahua Chen and two anonymous referees are gratefully acknowledged.

REFERENCES

- J. Chen & R. R. Sitter (1999). A pseudo empirical likelihood approach to the effective use of auxiliary information in complex surveys. *Statistica Sinica*, 9, 385–406.
- J. Chen, R. R. Sitter & C. Wu (2002). Using empirical likelihood methods to obtain range restricted weights in regression estimators for surveys. *Biometrika*, 89, 230–237.
- J.-C. Deville & C. E. Särndal (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association*, 87, 376–382.
- H. O. Hartley & J. N. K. Rao (1968). A new estimation theory for sample surveys. *Biometrika*, 55, 547–557.
- C. T. Isaki & W. A. Fuller (1982). Survey design under the regression superpopulation model. *Journal of the American Statistical Association*, 77, 89–96.
- A. B. Owen (1990). Empirical likelihood ratio confidence regions. *The Annals of Statistics*, 18, 90–120.
- A. B. Owen (2001). *Empirical Likelihood*. Chapman & Hall/CRC, New York.
- R. H. Renssen & N. J. Nieuwenbroek (1997). Aligning estimates for common variables in two or more sample surveys. *Journal of the American Statistical Association*, 92, 368–374.
- K. D. Zieschang (1990). Sample weighting methods and estimation of totals in the consumer expenditure survey. *Journal of the American Statistical Association*, 85, 986–1001.

Received 20 November 2002

Accepted 22 October 2003

Changbao WU: cbwu@uwaterloo.ca

Department of Statistics and Actuarial Science

University of Waterloo, Waterloo, Ontario, Canada N2L 3G1