

Methods used for Small Domain Estimation of Census Net Undercoverage in the 2001 Canadian Census

Peter Dick and Yong You

Peter Dick, Social Survey Methods Division, Statistics Canada, Ottawa, Canada. E-mail: Peter.Dick@statcan.ca.
Yong You, Household Survey Methods Division, Statistics Canada, Ottawa, Canada. E-mail: yongyou@statcan.ca.

Abstract

Since the 1991 the population estimates have used the Census counts adjusted for the net Census undercoverage. The population estimates require the net missed persons by single year of age for each sex for all provinces and territories are required. While the coverage studies can provide reliable estimates of missed persons for each province and territory, the sample size is not large enough to provide reliable estimates at the detail needed. A mixture of procedures is used to produce the estimates. Direct survey estimates for each province and territory are used to create one margin. A spline smoothing method is used to produce the national estimates of age and sex needed for the other margin. An Empirical Bayes regression model creates the estimates for broad age groups within a province. A synthetic model then generates the detailed single year of age estimates. Finally, a calibration procedure is used to ensure the detailed estimates are consistent with the fixed marginal totals. Of some concern is the measure of the quality of these estimates. A Mean Square Error (MSE) is produced for the Empirical Bayes regression model but this does not take into account all the adjustments made to the model. This paper will review and compare procedures for estimating the MSE for this small domain estimation problem.

KEY WORDS: Benchmarking, Calibration, Census undercoverage, Empirical Bayes, Hierarchical Bayes, Sampling variance, Synthetic estimation, Population estimation.

1. Introduction

The Census of Canada was conducted on May 15, 2001. One objective is to provide the Population Estimates Program with accurate baseline counts of the number of persons by age and sex for specified geographic areas. The count of persons includes usual residents, immigrants and non-permanent residents; excluded are all foreign visitors and non Canadian residents without a permit. Unfortunately, not all persons are correctly enumerated by the Census. Two errors that occur are undercoverage - exclusion of eligible persons - and overcoverage - erroneous inclusion of persons.

The main coverage vehicle used by Statistics Canada is the Reverse Record Check (RRC). This sample survey, with a sample size of 60,000 persons, estimates the net number of persons missed by the Census. This estimate is the combined total of the two types of coverage errors, the gross number of persons missed by the Census and the gross number of persons erroneously included in the final Census count. Once these estimates are adjusted for the coverage errors of persons living in collective dwellings, the final net number of people missed by the Census can be produced. The RRC sample size produces reliable direct estimates for large areas, such as provinces, and for large domains, such as broad age - sex combinations at the national level. However, the Population Estimates Program requires estimates of missed persons for single year of age for both sexes for each province and territory - over 2,000 estimates. Clearly the direct survey estimate would result in estimates having either unacceptably high standard errors due to insufficient sample in the small domain or having no estimate at all due to no sample in the domain. In addition, estimates have to be produced for the 288 Census Divisions and 4 different types of marital status. Altogether over 2.5 million estimates have to be created.

The current methodology used to generate these estimates has essentially been in place since 1991. One component of the procedure is to use the basic area model, as in Fay – Herriot (1979). However some modifications have been made to this basic model that needs to be evaluated. Specifically, the usual basic area model assumes that the sampling variances are known. The Census undercoverage model has to smooth the observed sampling variances before they can be used in the model. The final MSE does not take into account this estimation so clearly this approach has underestimated the uncertainty. Another drawback to the current methodology concerns the constraints that are imposed on the final estimates. Again the impact of this approach is to underestimate the MSE. The proper evaluation and impact of these two approaches is addressed in this paper. The chosen method is to adjust the model fit it into a Hierarchical Bayes framework. With this approach we can use the machinery developed over the last 10 years for evaluating this Hierarchical Bayes (HB) model and observe if the measures of uncertainty are comparable.

An advantage of the Hierarchical Bayes approach is that it is relatively straightforward and the inferences about the level parameters are “exact” unlike the Empirical Best Linear Unbiased Prediction (EBLUP) approach. The HB approach will automatically take into account the uncertainties associated with unknown parameters. However, it does require the specification of prior distributions. Fortunately the Census provides a case in which specifying the model is, again, relatively straightforward. The main purpose of this paper is to see how well the established method (EBLUP) compares with the more comprehensive HB approach.

Section 2 presents an overview of the methods that have been used to produce these estimates over the last 3 Census periods. This section gives a brief overview of the entire methodology but concentrates on the Empirical Best Linear Unbiased Prediction component. Section 3 presents the hierarchical Bayes models including a basic area level model and the modifications to the basic model for unknown sampling variances and constraining the final estimates.

2. Overview of Established Method

Mixtures of procedures are used to produce the required estimates of net missed persons. Direct survey estimates for each province and territory are used to create one margin. A cubic spline method is used to produce the national estimates of age and sex needed for the other margin. An Empirical Bayes regression model creates the estimates for broad age groups within a province. A synthetic model then generates the detailed single year of age estimates. Finally, a raking ratio procedure is used to ensure the detailed estimates are consistent with the fixed marginal totals. Dick (2001) gives a complete description of the approach used in 1996. Overall a summary giving with the background information of these methods can be found in Rao (2003).

(a) Marginal Totals

Two marginal totals are required: the provincial estimate of net missed persons and the national total of net missed persons by single year of age for each sex. The provincial and territorial estimates of net missed persons can be found in the official population release on coverage studies. These are assumed to be correct and all further estimates will respect these totals.

The direct estimates for national age and sex cannot be used without some smoothing. From the data, it appears that from age 0 to about age 15 there is very little pattern. Starting in the late teens until about age 30, there is a sharp increase in undercoverage. After age 30, there is gradual decline until about age 60 at which point the variability of the data becomes quite large. There are a large variety of methods that can be used to smooth the direct estimates. In the next section a smoothing method based on a generalized linear model is introduced but in this case it was felt that describing the model in a few parameters would be very difficult. Instead a nonparametric small domain model was introduced. This approach implicitly assumes that consecutive age groups have relatively similar undercoverage rates and change between ages follows a smooth function.

The model assumes that the true undercoverage rates - defined as the ratio of net missed persons over the total population - are described by a smooth function of age $f(a)$. Suppose we write the observed rate as equal to the smooth function at an age a plus a random error, or $r_a = f(a) + \varepsilon_a$. For this model we will assume that the function $f(a)$ is continuous - meaning that undercoverage rates move smoothly between consecutive ages. For reasons that will be shown, the model also assumes that the function has first and second derivatives and the error is assumed to be independent for different ages. More details on this model can be found in Ramsey (2000).

The RRC publishes the undercoverage rates by selected broad age groups: 0 to 19, 20 to 29, 30 to 44, 45 to 69 and 70 and over. The smoothing spline will produce different estimates for these groups. A common procedure is to ensure those small domain estimates are in agreement with higher level aggregated data is to calibrate the estimates to these totals. The spline estimates can be calibrated for bench marking quarterly or monthly time series to annual totals.

(b) Small Domain Estimation

The detailed estimates of missed persons by single of age and sex within each province is handled in a two step procedure. First, an Empirical Bayes regression model is used with broad age groups, and then a synthetic estimate is created for the single years of age within the broad age groups.

The objective of modelling the small domain estimates is to produce a series of estimates with a smaller Mean Square Error than the direct estimate. However, as opposed to the direct survey estimate which is design unbiased, the modelling approach will introduce a bias for each estimate. Thus modelling the small domain estimates implies that a trade off is required between reducing the variance of each estimate and the bias introduced through the modelling process. One approach to ensuring that the more reliable direct survey estimates are utilized is to introduce an Empirical Bayes model similar to Fay and Herriot (1979). This procedure creates an estimate that is a combination of a model estimate and the direct survey estimate weighted by their respective variances. It is an Empirical Bayes estimate instead of a Bayes estimate because underlying parameters are first estimated, then these estimated parameters are considered known in later calculations. Note that since the individual sampling variances are used in the estimation, a more precise direct estimate would contribute much more to the final Empirical Bayes estimate than a similar estimate with low precision. This ensures that the model does not dominate estimates that are already considered reliable. It is also possible to approach this estimation problem through a Hierarchical Bayes methodology. The approach is discussed in Section 3.

The Empirical Bayes regression will use the direct survey estimates of adjustment factors. If we write the Census count for age group j in province i as C_{ij} and the net missed persons - the difference between the undercount and the over count in the same domain - by $M_{ij} = U_{ij} - O_{ij}$ then we can define the net undercoverage rate to be $R_{ij} = M_{ij} / (M_{ij} + C_{ij})^{-1}$ and the adjustment factor - the ratio of true population to Census population - as $y_{ij} = (M_{ij} + C_{ij}) / C_{ij}$. Adjustment factors are used in the modelling since the true population in an area is simply the product of the Census count and the adjustment factor. Later, when synthetic estimates are required this property will be found to be very useful. Note that adjustment factors are related to the net undercoverage rates through $y_{ij} = (1 - R_{ij})^{-1}$.

The empirical Bayes model assumes that a two stage model. First, a sample model describes the basic relationship between the observed adjustment factors and the true adjustment factors. Secondly, a regression model describes how the true adjustment factors vary with a set of underlying variables. If we let θ_i be the true adjustment factor in a domain, then the model for the first stage or the sample survey model can be written as (dropping the subscript j)

$$y_i = \theta_i + \varepsilon_i,$$

where we assume the random error term is $\varepsilon_i \approx N(0, \sigma_i^2)$ and the variance is known. The second stage model relates the true adjustment factor to a set of underlying explanatory variables. If we let X_i be the set of variables then we can write the second stage model as

$$\theta_i = X_i \beta + \xi_i,$$

where the model term is assumed to be $N(0, \sigma_v^2)$. We also assume that the model errors are independent of the sampling errors and there is zero covariance between the two errors (i.e. $Cov(\varepsilon_i, \xi_i) = 0$). When we combine the models together we can estimate the expectation of the adjustment factors given the observed adjustment factors by

$$E(\hat{\theta}_i | y_i) = \hat{\gamma}_i y_i + (1 - \hat{\gamma}_i) X_i \hat{\beta}$$

where the shrinkage factor $\hat{\gamma}_i = \hat{\sigma}_v^2 / (\hat{\sigma}_i^2 + \hat{\sigma}_v^2)^{-1}$ determines the weights given to the direct estimate and the model estimate. Note that if the sampling error, σ_i^2 , is very small - implying little error associated with the direct estimate and consequently $\gamma_i \cong 1$ - the Empirical Bayes model will place more weight on the direct estimate. However, if the sampling error is very large - implying the direct estimate has a large amount of uncertainty attached to it and consequently $\gamma_i \cong 0$ - then the Empirical Bayes estimate places more weight on the model estimate. It is this property that makes this model so attractive for small domain estimation. The RRC sample sizes vary, by province, from about 3,000 persons in PEI to over 11,000 in Ontario. Consequently we would like to preserve the reliable direct estimates produced in Ontario and, in turn, rely more on the model estimate in PEI. This Empirical Bayes model allows for this.

One requirement for this model is that the sampling errors are assumed known. This adds a layer of difficulty to using this model in this situation. The sample variances for the direct survey estimates are estimated by a replicate methodology. The 5 replicates will produce unbiased but highly variable estimates of the sampling variance. If the direct estimates for the sampling variances are used in the model then the final estimates from the Empirical Bayes model are unstable. Thus, before using the model, the sampling variances must be stabilized. We will use $s_i^2 = \sigma_i^2$ for the estimate of the sampling variance.

One approach to estimating the sampling variance is to use the sample design of the RRC. Using the fact, that the sample was selected proportional to the population size then it is shown in Dick (1995) that to a reasonable approximation the log of the variance of missed persons can be written as $\log(V(M_i)) = \eta + \Gamma C_i + \zeta_i$, where $V(M_i)$ is the variance of the missed persons. Plots of the residuals do not seem to indicate any problems, so the straight line is taken as the estimate of the sampling variance. These estimates are assumed to be without any error and used in the Empirical Bayes regression model. See Dick (2001) for more details on evaluating this procedure.

The regression coefficients are estimated using a Weighted Least Squares method while the model variance is estimated by a Restricted Maximum Likelihood approach. An iterative scheme is used in SAS/IML which first estimates the regression parameters and then uses these estimates of the regression parameters to estimate the model variance. The updated model variance is then used in the next cycle for the weighted least squares. Convergence is fairly rapid with only about 6 or 7 iterations being needed. More details can be found in Dick (1995).

The final requirement of the Empirical Bayes model is produce an estimate of the Mean Square Error (MSE) of the estimates. If all the components were known, then the MSE would simply be the Bayes estimate of the posterior variance or $MSE(\hat{\theta}_i) = g_{1i} = \hat{\gamma}_i s_i^2$. Since $0 \leq \gamma_i \leq 1$ then, when the components are known, we are guaranteed the MSE of the Empirical Bayes model will be less - sometimes much less - than the original variance. However the components are not known but need to be estimated. Consequently, this term only gives a first (underestimated) approximation to the true MSE.

Prasad and Rao (1990) provide an approximation to the true MSE of the Fay - Herriot model when the components must be estimated. The MSE approximation is given as $MSE(\hat{\theta}_i) = g_{1i} + g_{2i} + 2g_{3i}$, where $g_{2i} = \hat{\sigma}_v^2 (1 - \hat{\gamma}_i)^2 x_i' (\Sigma \hat{\gamma}_i x_i x_i')^{-1} x_i$, which estimates the uncertainty due to estimating the regression parameter β and $g_{3i} = s_i^4 (\hat{\sigma}_v^2 + s_i^2)^{-3} V(\hat{\sigma}_v^2)$, which estimates the uncertainty due to estimating the model variance σ_v^2 . Note that the sampling variance is assumed known so any uncertainty concerning this variance is not included in the final MSE. This is an obvious short coming of this approximation.

Recently, Wang (2000), Wang and Fuller (2003), Rivest and Vandal (2002) and Datta, Rao and Smith (2002) have modified the basic Prasad - Rao estimator of the MSE to include terms that account for the unknown sampling variance and the estimation method of the model variance σ_v^2 . The impacts of these new approximations are discussed in You and Chapman (2003). They note that when the sample sizes are large the EBLUP methods perform well.

The requirements of the Population Estimates Program require that the estimate of missed persons be produced for single year of age - for each sex - for each province. The Empirical Bayes estimates, discussed above, can produce estimates for broad age groups - usually 4 or 5 ages - but it cannot work effectively with more age groups due to a lack of sample. Too many domains would produce an estimate of zero since no sample was in the domain. Clearly, to meet the requirements of the program, a synthetic estimate must be introduced.

Suppose, from the Empirical Bayes model, we have an estimated adjustment factor of y_{pk} for some broad age - sex group k within province p. The adjustment factors allow for an easy synthetic estimate of missed persons for any single year age, a, with the k-th age group by using $M_{pka} = C_{pka} (y_{pk} - 1)$. Note the assumption that for all the age groups with the k-th group, the same adjustment factor y_{pk} applies. This assumption allows for easy use of the adjustment factors but, in fact, is in contradiction to the earlier assumption concerning the undercoverage rates for single year of age. This implies that further refinements are needed before the estimate of missed persons can be used.

(c) Consistency Adjustment

The small domain estimation will no longer be consistent with the marginal totals discussed in Section 2(a). Hence a raking ratio adjustment is used on the small domain estimates to ensure consistency with both the provincial totals and the national age - sex totals. This procedure organizes the estimate of missed persons into a matrix with the single year of age estimates as the row and the province estimates as the columns. The fixed marginal totals are then used for the single year age at the national levels and for the provincial totals. The Empirical Bayes estimate of the adjustment factors are then used to generate the synthetic estimates of missed persons for each province. These estimates are then alternatively adjusted so that

they sum to the row and column totals. Convergence is usually reached in about 3 iterations. Complete details can be found in Dick (1995). These estimates are then used as the small domain estimate of missed persons. By adding them to the Census counts, the Population Estimates Program can create a baseline population for generating the population estimates. Further estimates may be required, such as for Census Divisions or for marital status. These are produced in a synthetic fashion using an appropriate adjustment factor as described in Section 2(b).

3. Hierarchical Bayes Model-Based Estimation

3.1 General Area Level Models

Let y_i denote the direct survey estimator of the i -th small area parameter of interest θ_i . Following You and Rao (2002), we consider a sampling model for y_i : $y_i = \theta_i + \varepsilon_i$, $i = 1, \dots, m$, with $E(\varepsilon_i | \theta_i) = 0$, that is, the direct survey estimator y_i is design-unbiased for the small area parameter θ_i . The sampling variance of y_i is $V(\varepsilon_i | \theta_i) = \sigma_i^2$. The sampling variance is usually assumed to be known in the model, but it may depend on the unknown parameter θ_i (You and Rao, 2002). The unknown parameter θ_i is assumed to be related to area level auxiliary variable x_i through a linking function g with random area effects v_i as $g(\theta_i) = x_i' \beta + v_i$, $i = 1, \dots, m$, where β is a vector of unknown regression parameters, and the v_i 's are uncorrelated with $E(v_i) = 0$ and $V(v_i) = \sigma_v^2$, where σ_v^2 is unknown. Normality of v_i is also assumed. The sampling model and the linking model are unmatched in the sense that they cannot be combined directly to produce a linear mixed effects model for small area estimation if the linking function g is a non-linear function (You and Rao, 2002).

3.2 Fay-Herriot model under HB framework

The Fay-Herriot model is a special case of the general model. In the Fay-Herriot model, the linking function $g(\theta_i) = \theta_i$ and the sampling variance σ_i^2 is replaced by a smoothed estimator $\tilde{\sigma}_i^2$ and then treated as known in the model. Under the HB framework, the Fay-Herriot model is given as (1) $y_i | \theta_i \sim \text{ind } N(\theta_i, \tilde{\sigma}_i^2)$, $i = 1, \dots, m$; (2) $\theta_i | \beta, \sigma_v^2 \sim \text{ind } N(x_i' \beta, \sigma_v^2)$, $i = 1, \dots, m$; (3) Priors: $\pi(\beta) \propto 1$, $\pi(\sigma_v^2) \sim IG(a_0, b_0)$. Inference based on the Gibbs sampling approach can be found in You and Rao (2002).

3.3 Benchmarked HB Estimators

You, Rao and Dick (2002) constructed benchmarked HB estimators for small areas. Let $\hat{\theta}_i^{HB}$ denote the HB estimator of θ_i and $\hat{V}(\theta_i)$ the posterior variance of θ_i . Let $\hat{\theta}_i^{BHB}$ denote the benchmarked HB (BHB) estimator of θ_i such that $\hat{\theta}_i^{BHB}$ is a function of the HB estimators $\hat{\theta}_i^{HB}$, $i = 1, \dots, m$, i.e., $\hat{\theta}_i^{BHB} = f(\hat{\theta}_1^{HB}, \dots, \hat{\theta}_m^{HB})$ for some function $f(\cdot)$, and satisfies the benchmark property: $\sum_{i=1}^m \hat{\theta}_i^{BHB} = \sum_{i=1}^m y_i$. For example, a ratio BHB (RBHB) estimator can be obtained as $\hat{\theta}_i^{RBHB} = \hat{\theta}_i^{HB} (\sum_{k=1}^m y_k) / (\sum_{k=1}^m \hat{\theta}_k^{HB})$. To obtain a measure of variability associated with the BHB estimator $\hat{\theta}_i^{BHB}$, we use the posterior mean squared error (PMSE), $\text{PMSE}(\hat{\theta}_i^{BHB}) = E[(\hat{\theta}_i^{BHB} - \theta_i)^2 | y]$, which is similar to the posterior variance associated with the HB estimator $\hat{\theta}_i^{HB}$. It can be shown (You, Rao and Dick, 2002) that the PMSE of $\hat{\theta}_i^{BHB}$ is given by $\text{PMSE}(\hat{\theta}_i^{BHB}) = (\hat{\theta}_i^{BHB} - \hat{\theta}_i^{HB})^2 + V(\theta_i | y)$. Thus the PMSE of $\hat{\theta}_i^{BHB}$ is simply the sum of the posterior variance $V(\theta_i | y)$ and a bias correction term $(\hat{\theta}_i^{BHB} - \hat{\theta}_i^{HB})^2$. The PMSE is readily obtained from the posterior variance and the estimators $\hat{\theta}_i^{HB}$ and $\hat{\theta}_i^{BHB}$.

3.4 Unknown sampling variances

The Fay-Herriot model assumes that the sampling variances σ_i^2 are known in the model. This is a very strong assumption. Usually a smoothed estimator of σ_i^2 is used in the model and then treated as known. In practice, the sampling variances

σ_i^2 are usually unknown and are estimated by unbiased estimators s_i^2 . The estimators s_i^2 are independent of the direct survey estimators y_i . Following Wang (2000), Rivest and Vandal (2002) and Wang and Fuller (2003), we also assume that $d_i s_i^2 \sim \sigma_i^2 \chi_{d_i}^2$, where $d_i = n_i - 1$ and n_i is the sample size for the i -th area. For example, suppose we have n_i observations from small area i and these observations are iid $N(\mu_i, \sigma_i^2)$. Let y_i be the sample mean of the n_i observations. Then $y_i \sim N(\mu_i, \sigma_i^2)$ and $\sigma_i^2 = \sigma^2 / n_i$. Then we can obtain an estimator of σ_i^2 as $s_i^2 = s^2 / n_i$, where s^2 is the sample variance of the n_i observations. Also y_i and s_i^2 are independent and $(n_i - 1)s_i^2 \sim \sigma_i^2 \chi_{n_i-1}^2$. We now present the Fay-Herriot model with the estimated sampling variances s_i^2 in a HB framework as follows: (1) $y_i | \theta_i, \sigma_i^2 \sim \text{ind } N(\theta_i, \sigma_i^2)$, $i = 1, \dots, m$; (2) $d_i s_i^2 | \sigma_i^2 \sim \text{ind } \sigma_i^2 \chi_{d_i}^2$, $d_i = n_i - 1$, $i = 1, \dots, m$; (3) $\theta_i | \beta, \sigma_v^2 \sim \text{ind } N(x_i' \beta, \sigma_v^2)$, $i = 1, \dots, m$; (4) Priors for the parameters: $\pi(\beta) \propto 1$, $\pi(\sigma_i^2) \sim IG(a_i, b_i)$, $i = 1, \dots, m$, $\pi(\sigma_v^2) \sim IG(a_0, b_0)$, where a_i, b_i ($0 \leq i \leq m$) are chosen to be very small known constants to reflect vague knowledge on σ_i^2 and σ_v^2 . IG denotes the inverse gamma distribution. For a complete HB inference, the Gibbs sampling method will be used. The full conditional distributions for the Gibbs sampler are given in You and Chapman (2003). The Rao-Blackwellized estimator of the posterior mean and posterior variances can be easily obtained from the full conditional distributions (You and Chapman, 2003).

References

- Datta, G. S., Rao, J.N.K. and Smith, D.D. (2002) On measures of uncertainty of small area estimators in the Fay – Herriot model. *Technical Report, University of Georgia, Athens*.
- Dick, J.P. (2001) Small domain estimation of missed persons in the 2001 Census. *Proceedings of the Survey Methods Section, Statistical Society of Canada* 37 – 46.
- Dick, J.P. (1995) Modelling net undercoverage in the 1991 Canadian Census. *Survey Methodology* 21: 44 – 55.
- Fay, R.E. and Herriot, R. A. (1979) Estimation of income for small places: An application of James-Stein procedures to census data. *Journal of the American Statistical Association* 74: 268 – 277.
- Prasad, N.G. N. and Rao, J.N.K. (1990) The estimation of the mean squared error of small area estimators. *Journal of the American Statistical Association* 85: 163 – 171.
- Rao, J.N.K. (2003) *Small Area Estimation*. John Wiley and Sons, New York
- Rivest, L.P. and Vandel, N. (2002) Mean squared error estimation for small areas when the small area variances are estimated. *Proceedings of the International Conference on Recent Advances in Survey Sampling, July 10 -13, 2002. Ottawa, Canada*
- Wang, J. (2000) Topics in Small Area Estimation with Applications to the National Resources Inventory, Ph.D. dissertation Iowa State University.
- Wang, J. and Fuller, W.A. (2003) The mean square error of small area predictors constructed with estimated area variances. *To appear Journal of the American Statistical Association*.
- You, Y. and Chapman, B. (2003) Small area estimation using area level models and estimated sampling variances. Statistics Canada Methodology Branch Working Paper (Draft).
- You, Y. and Rao, J.N.K. (2002) Small area estimation using unmatched sampling and linking models. *The Canadian Journal of Statistics* 30: 3 – 15.
- You, Y., Rao, J.N.K. and Dick, J.P. (2002) Benchmarking hierarchical Bayes small area estimators with application in census undercoverage estimation. *Proceedings of the Survey Methods Section 2002, Statistical Society of Canada* 81 - 86.