

AN APPLICATION OF REGRESSION AND CALIBRATION ESTIMATION TO POST-STRATIFICATION IN A HOUSEHOLD SURVEY

Bodhini Jayasuriya, Richard Valliant, U.S. Bureau of Labor Statistics
Richard Valliant, 2 Massachusetts Av., N.E., Rm 4915, Washington, D.C. 20212

Key Words: General regression estimator, Principal person method, Replication variance

ABSTRACT

This paper empirically compares three estimation methods—regression, calibration, and principal person—used in a household survey for post-stratification. Post-stratification is important in many household surveys to adjust for nonresponse and the population undercount that results from frame deficiencies. The correction for population undercoverage is usually achieved by adjusting estimated people counts in each post-stratum to equal the corresponding population control counts typically available from an external source such as a census. We will compare estimated means from the three methods and their estimated standard errors for a number of expenditures from the Consumer Expenditure Survey sponsored by the Bureau of Labor Statistics in an attempt at understanding how each estimation method accomplishes this step in post-stratification.

1. INTRODUCTION

In large household surveys, post-stratification is a means of reducing mean square errors by adjusting for differential response rates among population subgroups and frame deficiencies that often result in undercoverage of the target population. In general, the population is subdivided into groups (post-strata) at the estimation stage based on information that affect the response variables. The estimator is constructed in such a way that the estimated total number of individuals falling into each post-stratum is equal to the true population count. Post-stratum population counts are typically available from an external census for numbers of persons but not for numbers of households. If household estimates are needed, a single weight must be assigned to each household while using the person counts for post-stratification. Regression estimators of totals or means accomplish this by using person counts in each household's auxiliary data. Calibration estimation, with a least-squares distance function, is closely related to regression estimation but

possibility that each person in a household may have a different weight. The weight associated with the person is then assigned to the household. This method is difficult to analyze theoretically. The regression estimator discussed in this paper, while easily adjusted to the population under count, automatically produces a household weight that is not based on any particular person or its members. Lemaître and Dufour (1987) and Statistics Canada's use of the regression estimator are of great regard.

There are a growing number of precedents for the use of regression estimators in surveys both in the theoretical literature and in actual survey practice. Statistics Canada has incorporated the general regression estimator into its generalized estimation system (GES) software that is used in many of its surveys. Fuller, Loughin and Rao (1993) discuss an application to the USDA National Food Consumption Survey. One of the attractions of regression estimation is that many of the techniques in surveys including the post-stratification estimator mentioned above are special cases of regression estimators. It also more flexibly incorporates auxiliary information than other more common methods. Other works on regression estimation and post-stratification include Bethlehem and Keller (1987), Casady and Valliant (1987), Deville and Särndal (1992), Deville, Särndal, and Tillman (1993), and Zieschang (1990).

In this study we compare the regression estimator with the PP estimator currently in use at the Bureau of Labor Statistics (BLS). The ordinary least-squares regression estimator has the disadvantage that it can produce nonpositive weights. A number of ways are suggested in the literature on how to overcome this problem. The most flexible is the calibration method introduced by Deville and Särndal (1992) which can remove any nonpositive weights as well as control extreme weights. Calibration estimators produced by these new methods are also compared to the original regression estimator and the PP estimator.

In Section 2, the three different estimation methods are presented. Section 3 is an application of these methods to the Consumer Expenditure (CE) Survey at the household level in the same setting as in Zieschang (1990). We conclude with

2. REGRESSION, CALIBRATION, AND PRINCIPAL PERSON ESTIMATION

First, we give a brief introduction to the regression estimator. A sample s of size n is selected from a finite population U of size N . Let the probability of selection of the i^{th} unit be π_i . The sample could be two-stage and the unit could be either the primary sampling unit or the secondary sampling unit. There is no need here to complicate the notation with explicit subscripts for the different stages of sampling. Let the variable of interest be denoted by y and suppose that its value at the i^{th} unit, y_i , is observed for each $i \in s$. Assume the existence of K auxiliary variables x_1, x_2, \dots, x_K whose values at each $i \in s$ are available. Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iK})'$, for each $i \in U$, where x_{ik} denotes the value of the variable x_k at unit i . Let $\mathbf{X} = (X_1, \dots, X_K)'$ denote the K -dimensional vector of known population totals of the variables x_1, x_2, \dots, x_K . The regression estimator is then motivated by the working model ξ :

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_K x_{iK} + \varepsilon_i \quad (2.1)$$

or $i = 1, \dots, N$. Here, β_1, \dots, β_K are unknown model parameters. The ε_i are random errors with $E_{\xi}(\varepsilon_i) = 0$ and $\text{var}_{\xi}(\varepsilon_i) = \sigma_i^2$ for $i = 1, \dots, N$. The term "working model" is used to emphasize the fact that the model is likely to be wrong to some degree. In the CE, y_i might be the total food expenditures by the consumer unit (CU) and the x_{ik} 's might be various CU characteristics like numbers of people of different ages, or CU income, that have an effect on the CU's expenditure on food. The variance of expenditures might be dependent on CU size so that having σ_i^2 proportional to the number of persons in the CU might be reasonable. Then, a linear regression estimator of the population total of y is defined to be

$$\hat{y}_R = \hat{y}_{\pi} + (\mathbf{X} - \hat{\mathbf{x}}_{\pi})' \hat{\beta} \quad (2.2)$$

where \hat{y}_{π} denotes the π -estimator (or Horvitz-Thompson estimator) of the population total of y , i.e., $\hat{y}_{\pi} = \sum_{i \in s} a_i y_i$,

where $a_i = 1/\pi_i$. Also, $\hat{\mathbf{x}}_{\pi} = (\hat{x}_{1\pi}, \dots, \hat{x}_{K\pi})'$ is the vector of π -estimators of the population totals of the variables x_1, x_2, \dots, x_K and

$$\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_K)' = \left[\sum_{i \in s} \frac{a_i \mathbf{x}_i \mathbf{x}_i'}{\sigma_i^2} \right]^{-1} \sum_{i \in s} \frac{a_i \mathbf{x}_i y_i}{\sigma_i^2}. \quad (2.3)$$

Even if model (2.1) fails to some degree, \hat{y}_R will still have reasonable design-based properties because, even though

The regression estimator \hat{y}_R can also be expressed as a weighted sum of the sample y_i 's with i th weight,

$$w_i = a_i \left[1 + (\mathbf{X} - \hat{\mathbf{x}}_{\pi})' \left(\sum_{i \in s} \frac{a_i \mathbf{x}_i \mathbf{x}_i'}{\sigma_i^2} \right)^{-1} \frac{\mathbf{x}_i}{\sigma_i^2} \right].$$

From (2.4) it is easily seen that the known population totals are exactly reproduced for the auxiliary variables.

The estimator of β in (2.3) does not account for correlation among the errors in model (2.1). In many populations, units that are geographically near each other, e.g., CU's in the same neighborhood, may be correlated. Using a full covariance matrix \mathbf{V} may be more efficient than using a diagonal matrix. Using a full covariance matrix may be more optimal (e.g., see Casady and Valliant 1993, Valliant 1992). Though use of a full covariance matrix may lower the variance of $\hat{\beta}$, the elements of \mathbf{V} will depend on the particular y being studied, and estimation of \mathbf{V} is generally a nuisance. Consequently, it is interesting to consider the simple case of $\mathbf{V} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ which leads to (2.2). Note that when the design variance $\text{var}_p(\hat{y}_R)$ is estimated, it will be necessary to use a model that properly reflects clustering and other complexities.

The regression estimator has the disadvantage that the weights can be unreasonably large, small or even zero. The calibration estimators of Deville and Särndal (1992) introduced next, add constraints to restrict the size of the weights. Calibration estimators are formed by minimizing a distance F , between some initial weight a_i and final weight, subject to constraints. The constraints involve the available auxiliary variables thus incorporate them into the estimator. The regression estimator presented above is a special case of the calibration estimator in which F is defined to be the general distance (GLS) distance

$$F(w_i, a_i) = a_i c_i (w_i / a_i - 1)^2 / 2 \quad \text{for } i = 1, \dots, n, \quad (2.4)$$

where c_i is a known, positive weight (e.g., $c_i = \sigma_i^2$ or $c_i = 1$) associated with unit i , and w_i , the final weight. The total distance $\sum_{i \in s} F(w_i, a_i)$ is minimized subject to the constraints, $\sum_{i \in s} w_i \mathbf{x}_i = \mathbf{X}$. In this form, the weighted regression estimator of the population total of y (2.4) can be written as,

$$w_i = a_i g(c_i^{-1} \lambda' \mathbf{x}_i)$$

for $i = 1, \dots, n$ where

$$g(u) = 1 + u,$$

for $u \in \mathfrak{R}$ and λ is a Lagrange multiplier evaluated in the minimization process. The calibration weights can be unreasonably extreme values resulting in

chosen in such a way as to reflect the desired restrictions on the weights. Choosing $L > 0$ ensures that the weights are positive, and U is picked to be appropriately small to prohibit large weights. The calibration weights must be solved for iteratively; one easily programmed algorithm is given in Stukel and Boyer (1992).

In most household surveys, post-stratification serves primarily as an adjustment for undercoverage of the target population by the frame and the sample. In the U.S., there are no reliable population counts of households to use in post-stratification. Consequently, population counts of persons are used for the post-strata control totals. This disagreement in the unit of analysis (the household) and the unit of post-stratification (the person) when a household characteristic is of interest led to the development of the PP method that is used in the CE and Current Population surveys.

In the PP method described in Alexander (1987), a household begins the weighting process with a single base weight, a_i , that is then adjusted for nonresponse. The adjusted weight is assigned to each person in the household and the person weights are then further adjusted to force them to sum to known population controls of persons by age, race, and sex. This last adjustment can result in persons having different weights within the same household. The household is then assigned the weight of the person designated as the "principal person" in the household. This method has an element of arbitrariness and is difficult to analyze mathematically. The regression and calibration estimators can be formulated in such a way that population person controls are satisfied, all persons in a household retain the same weight, and no arbitrary choice among person weights is needed to assign a household weight.

3. AN APPLICATION

We compare the three estimators (i.e., regression, restricted calibration (with $L=.5$, $U=4$), and principal person) by an application to the estimated means and their estimated standard errors for a number of expenditures from the CE Survey sponsored by the Bureau of Labor Statistics.

The CE Survey gathers information on the spending patterns and living costs of the American consumers. There are two parts to the survey, a quarterly interview and a weekly diary survey. The Interview Survey collects detailed data on the types of expenditures which respondents can be expected to recall for a period of three

$n = 5156$ CU's were used. The CE Survey's unit of analysis is the consumer unit, an economic family household. A consumer unit (CU) consists of individuals in the household who share expenditures. Thus, there can be more than one CU in a household.

Five different sets of auxiliary variables were used. They were chosen by testing the adequacy of models for the selected expenditures with different combinations of the available auxiliary variables. The 56 post-stratification age/race/sex currently in use in the CE were used. The combinations of auxiliaries used to form the weights are given in Table 1. The number of variables in each model is given within parentheses on this information, weights (2.5) were computed as given in (2.6)—regwts—and (2.7)—calwts. For regression and restricted calibration weights, weights are equal to the adjusted base weight, i.e., $1/\pi_i$, with a nonresponse adjustment.

Table 1. Weights and their corresponding auxiliary variables. Number of cells are in parentheses.

Weight	Auxiliary Variables
regwts0	age/race/sex (56)
regwts1	inter., age/race/sex, region, urban×region (18)
regwts2	intercept, age/race/sex, region, urban×region, age of reference person, housing tenure, family income before taxes (24)
calwts0	age/race/sex (45)
calwts1	inter., age/race/sex, region, urban×region (18)
calwts2	intercept, age/race/sex, region, urban×region, age of reference person, housing tenure, family income before taxes (24)
calwts3	intercept, age/race/sex, region, urban×region, family income before taxes (truncated at \$500,000) (19)
calwts4	intercept, age/race/sex, region, urban×region, age of reference person, housing tenure (23)
PP	age/race/sex (56)

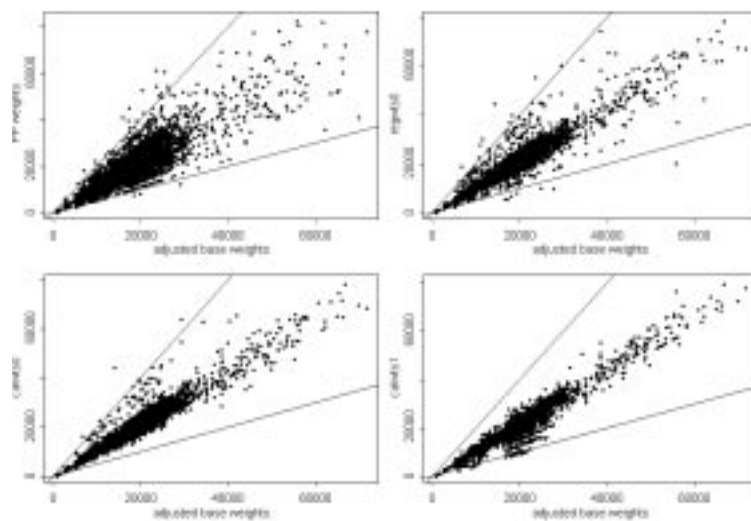
For this application, the population totals needed to evaluate $\mathbf{X} = (X_1, \dots, X_K)'$ were obtained mostly from 1990 Census figures projected to 1992 and the Population Reports published by the U.S. Bureau of Census.

3.1 Comparisons of Weights and Estimated Cell Counts

eplicate weights, nearly half the sets for each of regwts0, egwts1 and regwts2 had some negative weights though the maximum number of negative weights for any replicate was 3. The negatives are a potential cause of inflated standard errors, since the negative weights will be offset by large positive weights in order for the fixed population control totals to be met in every replicate. Calwts, which restrict the deviation from the base weights by choosing L and U appropriately, (in this instance, $L=0.5 > 0$) naturally did not produce any negative weights.

On examining scatter plots (not shown here) comparing some of the different weights to each other, the PP and regwts0, while being substantially different from each other, exhibited final weights that can be considerably different from the adjusted base weights. The adjustments can be either up or down. A less variable set of adjustments was apparent in regwts1, calwts0, and calwts1. Calwts1 and calwts4 were quite similar and both were close to regwts1. The two sets of weights that involve the quantitative variable family income before taxes, calwts2 and calwts3, were closely related. Some CU's had calwts2 values larger than 60,000 but had calwts0, calwts1, calwts4 < 30,000. These CU's all had family incomes before taxes of a quarter of a million dollars or more. Thus, the inclusion of that variable in the calibrations did have a substantial impact on some units. We did use a control only on the grand total income; having controls by income classes might have changed the weights on some of these cases.

Figure 1. Four sets of weights plotted against adjusted base weights. Reference lines correspond to $L=0.5$ and $J=2$.



indicate that the PP weights and regwts0 do not meet the restriction $a_i / 2 \leq w_i \leq 2a_i$.

Previous studies at BLS regarding regression estimation in the CE had concluded that the number of single person CU's was underestimated compared to the estimate produced by the PP method. We found minimal evidence of that phenomenon here, as indicated by the ratios shown in Table 2. It is a 10% ratio of the estimated number of CU's under the alternative procedures to that of the PP estimation procedure for a single person CU.

Table 2. Estimated counts in thousands of CU's by size for PP weights and ratios of other estimated counts to the PP weights estimates. Ratios greater than 1.0% and less than 0.98 are highlighted.

Weights	CU Size				
	1	2	3	4	5+
PP	28,784	30,680	15,409	15,068	9,993
regwts0	0.96	0.99	1.01	0.99	1.02
regwts1	1.00	1.00	1.02	0.98	0.99
regwts2	1.00	1.01	1.01	1.01	0.97
calwts0	0.96	0.99	1.00	1.00	1.02
calwts1	1.00	1.00	1.02	0.98	0.99
calwts2	0.99	1.01	1.01	1.00	0.97
calwts3	0.98	1.02	1.01	1.01	0.96
calwts4	1.01	1.00	1.01	0.99	0.99

A similar table constructed by Composition of CU that while regwts0 and calwts0 estimates are substantially different from PP for the category Or 1+ children, for single person CUs they were not.

3.2 Precision of Estimates from the Different Methods

Although comparison of weights is instructive, the different methods must ultimately be judged based on the estimated CU means and their precision. The standard errors of these estimators were computed via the method of balanced half sampling (BHS) using 44 replicates, which is currently implemented in the CE for the PP estimator. The BHS estimator is constructed to reflect the stratification and the clustering that is used in the CE. For the expenditure estimates from the CE Survey are available for various domains of interest, we computed the means and the standard errors for a few chosen domains. For each of these, the coefficient of variation was computed and then its ratio to the cv of the PP estimator was calculated.

expenditures, and for each of the following domains: Age of Reference Person, Region, Size of CU, Composition of Household, Household Tenure, and Race of Reference Person.

Table 3. Ratios to CE cv to cv's for the different weighting methods. The minimum ratio is highlighted in each row.

Expenditure	regwts			calwts		
	0	1	2	0	1	2
All Exp.	0.98	0.90	0.79	0.98	0.90	0.78
Shelter	0.93	0.85	0.75	0.93	0.85	0.74
Utilities	1.08	1.03	0.94	1.07	1.03	0.88
Furniture	1.08	1.21	3.52	1.06	1.21	2.58
Maj. ap.	1.08	1.06	1.04	1.06	1.08	1.09
All vehi.	0.90	0.89	0.98	0.91	0.89	0.98
Cars, (n)	0.95	0.91	1.01	0.96	0.91	1.02
Cars, (u)	0.98	0.94	0.96	0.97	0.94	0.97
Gasol.,	1.17	1.11	1.03	1.12	1.10	0.99
Health	1.05	0.97	0.86	1.07	0.97	0.85
Educat.	0.92	0.93	1.04	0.91	0.93	1.06
Contrib.	1.01	1.02	1.28	1.01	1.02	1.30
Pers.	1.00	0.97	1.64	1.01	0.98	1.24
Ins.						
Life, Ins.	1.08	1.02	1.53	1.08	0.98	1.38
Pensions	1.00	0.99	1.75	1.01	0.99	1.34

In addition, ratios for all CU's, i.e., the total across the domains, were computed for each expenditure and those for regwts 0, 1, 2 and calwts 0, 1, 2 are shown in Table 3. For All Expenditures, regwts2 and calwts2 with ratios of .79 and .78 provide substantial reduction in cv compared to PP. For less aggregated expenditures, regwts1 or calwts1 provide reasonably consistent improvements over PP without the losses incurred by some of the other weights or expenditures like Furniture, Personal insurance and Pensions, and its subcategory Pensions and social security.

A trellis plot (Cleveland 1993) of the cv and mean ratios for calwts0 and calwts1 by age of reference person is given in Figure 2. Calwts0 is pictured because it is the nearest calibration equivalent to the current method of post-stratification. Calwts1 appears to be the best of the alternatives we have examined in the sense of improving the All Expenditures estimates while providing consistent performance for individual expenditure groups. In each panel of the plot a vertical reference line is drawn at 1, the point of equality between the calibration results and those of the PP method. The lower tier in the plot presents

ratios tend to be less than 1, for most domains, and calwts1 is somewhat better than

Calwts2 and calwts3, which used family income taxes as one of the auxiliaries, had somewhat better performance for domains, sometimes making improvements over PP but occasionally showing losses. This is connected to the nature of the income variable which had a substantial number with negative and zero values. These CU's suggest the usefulness of this variable in predicting expenditures.

Taking all of the above into consideration, calwts1 and calwts4 can be deemed a clear improvement over the PP estimator. Calwts1 has the advantage of negative weights over regwts1. Since calwts4 requires auxiliary variables as opposed to calwts1's weights, we recommend calwts1 over all the other types of weights we have considered.

4. CONCLUSION AND FUTURE RESEARCH

The objective of this study was to investigate alternatives to the principal person method for household weights that did not depend on the presence of one single member of the household. Different weighting methods based on the regression estimation procedure were presented and their relative merits evaluated. Regression estimation incorporates the current survey stratification methods in which the weighted sum of persons in each post-stratum is forced to be equal to the independent census count of that number. This was accomplished via auxiliary variables that are incorporated into the regression model. It also automatically provides for each sample household a weight that does not depend on any single one of its members. In order to eliminate undesirable negative weights that can result from least-squares regression estimation, calibration weights were adapted to the present problem. The current estimation procedure has the flexibility to reduce the possible deviation of each final weight from its base weight while adhering to the properties discussed above. In particular, it allows the constraint of positive weights. Calibration weights are easily computed via matrix algebra using software like S-PlusTM.

Overall, the ordinary regression estimator and the calibration estimator both appeared to be an improvement over the Principal Person estimator in terms of reducing the coefficient of variation. For the future, the calibration estimator can be further refined by using the procedure of regression estimation to choose the auxiliary variables.

REFERENCES

- ALEXANDER, C. H. (1987). A class of methods for using person controls in household weighting. *Survey Methodology*, 13, 183-198.
- BETHLEHEM, J.G., and KELLER, W.J. (1987). Linear weighting of sample survey data. *Journal of Official Statistics*, 3, 141-153.
- CASADY, R J., and VALLIANT, R. (1993). Conditional properties of post-stratified estimators under normal theory. *Survey Methodology*, 19, 183-192.
- LEVELAND, W.S. (1993). *Visualizing Data*. Summit, New Jersey: Hobart Press.
- DEVILLE, J.C., and SÄRNDAL, C.E. (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association*, 87, 376-382.
- DEVILLE, J.C., SÄRNDAL, C.E., AND SAUTORY, O. (1993). Generalized Raking Procedures In Survey Sampling. *Journal of the American Statistical Association*, 88, 1013-1020.
- FULLER, W. A., LOUGHIN, M. M. and BAKER, H. D. (1993). Regression weighting for the 1987-1988 Nationwide Food Consumption Survey. Unreport submitted to the United States Department of Agriculture.
- LEMAÎTRE, G. and DUFOUR, J (1987). An improved method for weighting persons and families. *Methodology* 13, 199-207
- RAO, J.N.K. (1992). Estimating totals and distribution functions using auxiliary information at the design stage. Presented at the Workshop on Uses of Information in Surveys, Statistics Sweden.
- SÄRNDAL, C. E., SWENSSON, B. and WRETENBERG, K. (1992). *Model assisted survey sampling*. New York: Springer - Verlag.
- STUKEL, D. M. and BOYER, R. (1992). Calibration estimation: an application to the Canadian labour force survey. Working Paper, Statistics Canada : SS 009E, 1992
- ZIESCHANG, K. D. (1990). Sample weighting and estimation methods in the Consumer Expenditure Survey. *Journal of the American Statistical Association*, 85, 986-1001.

Figure 2. Ratios to CE of cv's and means for two weighting methods by age of reference person.