**UNITED NATIONS STATISTICAL COMMISSION and
ECONOMIC COMMISSION FOR EUROPE
CONFERENCE OF EUROPEAN STATISTICIANS**

**Work Session on Statistical Data Editing**
(Madrid, Spain, 20-22 October 2003)

Topic (i): Development and use of data editing quality indicators

**TITLE**

**Performance of Bootstrap Techniques with Imputed Survey Data**

Submitted by INE, Spain[1]

## I.     INTRODUCTION

1.      In this study we analyse the performance of the bootstrap, when estimating the variance and a confidence interval of a simple statistic in a survey with nonresponse and mean imputation. The relative bias and root mean square error of the variance estimator and also the coverage rate of the bootstrap confidence interval are calculated and compared with the ones obtained from the jackknife variance estimator in a previous work (Aparicio-Perez, F. and Lorca, D. 2002).

2.      The data are taken from the Industrial Business Survey and are exactly the same as in our previous work.

3.      In the following section, we explain the bootstrap methods that are used, next we present the simulation study and results and finally some conclusions are provided.

## II.     BACKGROUND

4.      Shao and Sitter (1996) proposed a bootstrap method for variance estimation with imputed survey data. This is the method that we are using. There are another modified bootstrap methods, like the ones proposed in Saigo, Shao and Sitter (2001). See also Lee, Rancourt and Särndal (2001) for a review of variance estimation methods under imputation.

5.      Assuming a rather general sampling design, the method consists of the following steps:

  a)  Let $Y_I = Y_R \cup Y_O$ where $Y_R = \{y_k : k \in A_R (respondents)\}$ and
      $Y_O = \{z_k : k \in A_O (non-respondents)\}$. $Y_O$ is obtained from $Y_R$ using some imputation technique.

  b)  Let $Y^* = \{y_i^* : i=1,\ldots n\}$ a simple random sample (bootstrap sample) drawn with replacement from $Y_I$, and $Y_I^* = Y_R^* \cup Y_O^*$ where $Y_R^* = \{y_k^* : k \in A_R^* (respondents)\}$/and

---

[1] Prepared by Félix Aparicio-Pérez and Dolores Lorca, (fapape@ine.es, mdlorca@ine.es).

$Y^*_O = \{y^*_k : k \in A^*_O (non-respondents)\}$ and $Y^*_O$ is obtained from $Y^*_R$ using the same imputation technique that was used in step a)

c) Obtain the bootstrap estimator $\hat{\boldsymbol{q}}^*_I = \hat{\boldsymbol{q}}(Y^*_I)$ of $\hat{\boldsymbol{q}}_I = \hat{\boldsymbol{q}}(Y_I)$, based on the imputed bootstrap data set $Y^*_I$.

d) Repeat steps b) and c) B times. The following bootstrap variance estimator for $\hat{\boldsymbol{q}}_I$ is calculated:

$$v_B(\hat{\boldsymbol{q}}_I) = \frac{1}{B} \sum_{b=1}^{B} (\hat{\boldsymbol{q}}^{*b}_I - \vec{\boldsymbol{q}}^*_I)^2$$

where

$$\vec{\boldsymbol{q}}^*_I = \frac{1}{B} \sum_{b=1}^{B} \hat{\boldsymbol{q}}^{*b}_I$$

6. We are also calculating the percentile and the bootstrap-t confidence intervals of the statistic. The former is computed as the empirical confidence interval obtained after sorting the B bootstrap statistics.

To build the bootstrap-t confidence intervals, we first compute the statistics $t^*_b = (\hat{\boldsymbol{q}}^{*b}_I - \vec{\boldsymbol{q}}^*_I) \big/ \boldsymbol{s}^*_b$ where

$\boldsymbol{s}^{*2}_b = v_B(\hat{\boldsymbol{q}}^{*b}_I)$, and then calculate $t^*_L = C\hat{D}F_t^{-1}(\boldsymbol{a})$, $t^*_U = C\hat{D}F_t^{-1}(1-\boldsymbol{a})$, where $C\hat{D}F_t^{-1}(x) = \#(t^*_b; t^*_b \le x, b=1....B)/B$. The bootstrap-t confidence interval is then given by $(\hat{\boldsymbol{q}}_I - t^*_U v_B(\hat{\boldsymbol{q}}_I), \hat{\boldsymbol{q}}_I - t^*_L v_B(\hat{\boldsymbol{q}}_I))$

## III. MONTECARLO STUDY AND RESULTS

7. From our population of N=16438 industrial businesses we draw simple random samples without replacement of sizes n=100, 500, 1000 and 5000. We simulate a loss of about 30 per cent of our data with a uniform mechanism. The number of replications is 50000 for each sample size. The number of bootstrap samples is 999 for each replication.

8. An additional simulation with 23000 replications was conducted only for small sample sizes (n=100) calculating the bootstrap-t confidence interval. We used 999 bootstrap samples in the first level and 50 bootstrap samples in the second level (variance estimation of the first level bootstrap statistics).

9. The analysis variable is the turnover of the businesses. We use mean imputation.

10. Within each replication, we compute the percentage relative bias, the relative root mean square error and the coverage of a nominal 95% bootstrap confidence interval.

**TABLE 1. Imputed variable: turnorver.**
Imputation method: mean imputation, 95% confidence intervals

| Sample size | BOOTSTRAP T | | | BOOTSTRAP Percentile | | | JACKKNIFE | | |
|---|---|---|---|---|---|---|---|---|---|
| | RB(%) | MSE | COVR(%) | RB(%) | MSE | COVR(%) | RB(%) | MSE | COVR(%) |
| 100 | 0.07 | 4.51 | 56.9 | -0.42 | 4.47 | 60.6 | 0.29 | 4.52 | 56.9 |
| 500 | | | | 1.00 | 2.00 | 68.3 | -0.90 | 1.96 | 65.8 |
| 1000 | | | | 4.24 | 1.45 | 75.6 | -2.12 | 1.35 | 72.6 |
| 5000 | | | | 27.15 | 0.76 | 91.9 | -12.0 | 0.51 | 87.3 |

## IV.    CONCLUSIONS

11.    The results show that the percentile bootstrap performs better than the jackknife for coverage rate of the confidence intervals and the reverse is true for mean square errors and bias of the variance estimators.

12.    There seems not to be any advantage in using the bootstrap-t confidence interval, in spite of its higher computational time, though only small sample simulations and small number of second level bootstrap samples are considered for this method.

## V.    REFERENCES

Aparicio-Pérez, F. and Lorca, D. (2002). Performance of jackknife variance estimation using several imputation methods. UNECE work session on Statistical data Editing. Helsinki.

Lee, H., Rancourt, E. and Särndal, C.E. (2001). Variance estimation from survey data under single imputation. In Survey Nonresponse. Eds. R.M. Groves, D.A. Dillman , J.L. Eltinge and R.J.A. Little, 315-328. Wiley. New York.

Saigo, Shao and Sitter (2001). A repeated half-sample bootstrap and balanced repeated replications for randomly imputed Data. Survey Methodology, vol 27,No. 2,pp.189-196

Shao and Sitter (1996) Bootstrap for imputed survey data. Journal of American Statistical Association,91,1278-1288.