# Rounding Non-integer Weights in Bootstrapping Non-iid Samples: actual problem or harmless practice?

Federico Andreis[1] and Fulvia Mecatti[2]

[1] DEMM - Dipartimento di Economia, Management e Metodi Quantitativi
Università degli Studi di Milano
Via Conservatorio 7, 20100 Milano, Italy
(e-mail: `federico.andreis@unimi.it`)
[2] Dipartimento di Sociologia e Ricerca Sociale
Università degli Studi di Milano-Bicocca
Via Bicocca degli Arcimboldi 8, U7, 20126 Milano, Italy
(e-mail: `fulvia.mecatti@unimib.it`)

ABSTRACT. The effect of the practice of rounding non-integer weights in bootstrapping complex samples form finite populations is investigated by means of an extended simulation study. The extent to which rounding can interfere with basic bootstrap principles as well as with the formal properties of the final bootstrap estimate are evaluated. Indications and recommendations on application of this method are discussed.

## 1 INTRODUCTION AND MOTIVATION

Bootstrap method is a popular tool for numerically assessing estimators accuracy, confidence intervals and p-values. In order to provide reliable results for survey sampling inference, feasible adaptations of the basic bootstrap algorithm accounting for the non-iid nature of data are required. Recent literature suggests the use of a weighting system in the resampling and/or estimation procedure (Antàl and Tillé (2011), Beaumont and Patak (2012), Ranalli and Mecatti (2012)). Integer weights guarantee analytical properties however in practical applications this is seldom the case. Two suggestions recur in the literature to deal with non-integer weights: randomization and systematical rounding. Randomization would require a further step added on top of the boostrap algorithm, thus affecting its computational efficiency. This step can be avoided by rounding non-integer weights to the nearest integer, according to some systematical rule (Chauvet (2007), Barbiero and Mecatti (2009)). Both solutions affect the bootstrap process as well as the final boostrap estimates to an unknown extent; moreover, they violate basic bootstrap principles such as the mimicking principle and the plug-in approach. In this work we concentrate on rounding for dealing with non integer weights as the computationally preferable option, our main aim being to produce empirical evidence of the extent of its effects.

## 2 METHODS

We focus on non-iid samples of size $n$ from a finite population $U = \{1, ..., k, ..., N\}$ selected without replacement under a random design where each population unit is assigned a spe-

cific probability $\pi_k$ to be included in the sample. We adopt the popular approach based on weighting each sampled unit by the inverse of its own inclusion probability $\pi_k^{-1}$ (Holmberg 1998) in order to produce an empirical population $U^*$ where to perform resampling. The pseudo-population $U^*$, usually named the *bootstrap population*, is intended to mimick the *parent population U* according to fundamental bootstrap principles such as the mimicking and the plug-in approach (Hall (1992)). According to the same principles, the resampling design should mimick the original sampling. The creation of the bootstrap population, the resampling and the final estimates produced by the bootstrap process depend on the weights $\pi_k^{-1}$, whether integer or to be rounded.

We have investigated the magnitude of the rounding effect on the bootstrap estimate of the variance of the Horvitz-Thompson (HT) estimator for the mean of a study variable through an extended simulation exercise.

A toy example that grants inclusion probabilities leading to integer weights was created and several scenarios were considered. Particularly gradual departures were generated from the 'perfect' integer-weights situation by inducing small to large deviations in the set of inclusion probabilities so that increasing degrees of rounding were simulated. Two sampling designs have been considered:

1. constant inclusion probabilities, $\pi_k = n/N, \forall k \in U$; and
2. unequal inclusion probabilities, for instance proportional to an auxiliary variable, $\pi_k \propto x_k$.

Notice that in case 1, rounding affects a single value, i.e. $n/N$ which is the same for all units in $U$, and a unique control total, i.e. the bootstrap population size $N^*$ possibily differing from the actual population size $N$. In case 2 the effect might be more severe, since rounding may occur for multiple values, namely form 1 to $n$ bootstrap weights. Moreover, in case 2 the rounding would affect two control totals, both the bootstrap population size $N^*$ and the total of the auxiliary variable $X^*$ which can depart from the corresponding actual population counts $N$ and $X = \sum_{k \in U} x_k$.

Numerical methods have been developed in order to find suitable set of inclusion probabilities leading to integer weights in order to build the toy example in both cases 1 and 2, the latter being clearly more complex so that requiring some analytical and computational effort. As for the departures from the ideal integer-weights situation when no rounding is needed, increasingly strong perturbations on the $\pi_k$ were induced: slightly varying the population size (case 1) or by means of additive stochastic noise from a Uniform distribution with suitable support (case 2).

## 3   SIMULATION DESIGN

Population values were generated from a Gamma distribution. In order to investigate scenarios with different variability, several combinations of the distribution parameters were chosen leading to $c.v. = 0.5, 1.0$ and $1.5$. A 5% sampling fraction was chosen in the integer weights case which has been kept within a tight range $(0.05 \pm 0.003)$ in the non-integer case. All the computations were carried out in *R* 2.15.2 code, and the packages *sampling* (Tillé and Matei (2012)) and *BiasedUrn* (Fog (2013)) were respectively employed to select unequal probability samples and to generate values from non-standard distributions as involved in this study.

The bootstrap methods compared, based on bootstrap population and plug-in approach, are Chao and Lo (1985) and Ranalli and Mecatti (2012) for case 1 (constant $\pi_k$) and Holmberg (1998) and Ranalli and Mecatti (2012) for case 2 (unequal $\pi_k$).

A set of Monte Carlo indicators was provided for each simulated scenario with the purpose of investigating the rounding problem under the following respects:

- the mimicking principle: by evaluating distances between the nominal and the post-rounding characteristics of the bootstrap algorithm, particularly on known population totals and size as compared to bootstrap populations counts;
- basic bootstrap algorithm properties: particularly the so called bootstrap unbiasedness as measured by Monte Carlo expectation over HT estimates computed on both the original sample and the colection of bootstrap samples;
- inferential properties of the final bootstrap estimates for the variance of the HT estimator, such as biasedness and stability as measured by Monte Carlo relative bias and relative root mean square error; and
- how good a simulation the bootstrap distribution is for the unknown estimator distribution in order to provide accurate confidence intervals, by comparing standard normal and bootstrap percentile confidence interval.

Simulation results give empirical ground for indications and recommendations about rounding in non-iid bootstrap algorithms either as an actual problem or as an harmless practice.

## REFERENCES

ANTÁL, E., TILLÉ, Y. (2011): A Direct Bootstrap Method for Complex Sampling Designs From a Finite Population. *Journal of the American Statistical Association*, *106*, 534–543.

BARBIERO, A., MECATTI, F. (2009): Bootstrap algorithms for variance estimation in PS sampling. In: P. Mantovan and P. Secchi, *Complex Data Modeling and Computationally Intensive Statistical Methods*. Springer-Verlag, Berlin.

BEAUMONT, J-F., PATAK, Z. (2012): On the Generalized Bootstrap for Sample Surveys with Special Attention to Poisson Sampling. *International Statistical Review*, *80-1*, 127–148.

CHAO, M.T., LO, A.Y. (1985): A bootstrap method for finite population. *Sankhya*, *47(A)*, 399–405.

CHAUVET, G. (2007): Méthodes de bootstrap en population finie. PhD Dissertation, Laboratoire de statistique d'enquêtes, CREST-ENSAI, Université de Rennes 2. Available at http://tel.archives-ouvertes.fr/docs/00/26/76/89/PDF/thesechauvet.pdf.

HALL, P. (1992): *The bootstrap and Edgeworth Expansion*. Springer-Verlag, New York.

HOLMBERG, A. (1998): A bootstrap approach to probability proportional to size sampling. In: *Proceedings of Section on Survery Research Methods*. American Statistical Association, 181–184.

R CORE TEAM (2013): R: A Language and Environment for Statistical Computing. *R Foundation for Statistical Computing, url: http://www.R-project.org*.

RANALLI, M.G., MECATTI, F. (2012): Comparing Recent Approaches For Bootstrapping Sample Survey Data: A First Step Towards A Unified Approach. In: *Proceedings of Section on Survery Research Methods*. American Statistical Association, 4088–4099.