Simplified variance estimation for multistage sampling - Application to longitudinal surveys

Guillaume Chauvet

Ensai-IRMAR

26/10/2018

- Context of this work
- Simplified variance estimation for two-stage sampling
- Bootstrap variance estimation
- Cross-sectional Estimation at Initial time
- $oldsymbol{5}$ Cross-sectional estimation at time t+1 : Panel and refreshment sample
- Related/further work

Context of this work

First work : Panel Politique de la Ville

The survey "Panel Politique de la Ville" (PPV) has been set up to study the life conditions for the inhabitants in the districts of the sensitive urban areas (ZUS). The four waves of survey from 2011 to 2014 aim at focusing on :

- residential mobility inside districts,
- the perception of public policies,
- the impact of public policies on recipients.

The original sample is selected through two-stage sampling (Couvert et al., 2016):

- drawing of a stratified sample of districts,
- inside districts, drawing of a sample of households by using a sampling frame built from the Census,
- inside households, all the individuals are surveyed.



Second work: Household Finance and Consumption Survey (HFCS)

The HCFS is conducted in a decentralised manner, though some common guidelines are provided by the HCF Network (HCFN) + coordination by the European Central bank. The methodology depends on the country: from 1 to 3 stages of sampling.

The HFCS provides detailed household-level data on different aspects of household balance sheets and associated economic and demographic variables. This includes income, employment and measures of consumption.

The priority of the HCFS is on cross-sectional estimations (i.e., at a given point in time).

The HFC Network (HFCN) also recommends the introduction of a panel component to measure individual changes over time.

Some common challenges

At initial time t, the sample is selected by means of a multistage sampling design. At further times, a split-panel strategy is used :

- tracking of a part of the units (individuals or households) surveyed at former time,
- selection of a refreshment sample of households and individuals to represent the newborns.

Purpose of my work:

- weighting for cross-sectional estimations,
- weighting for longitudinal estimations,
- assorted variance estimation :
 - linearization variance estimation for PPV,
 - bootstrap variance estimation for HCFS.

In this talk: focus on cross-sectional estimation + bootstrap



Simplified variance estimation for two-stage sampling

Two-stage sampling design

We are interested in a finite population U of units, say households. These households are grouped into larger Primary Sampling Units (PSUs), say municipalities. We note U_I the population of PSUs.

We are interested in the total

$$Y = \sum_{k \in U} y_k$$
 for some variable of interest y_k ,
 $= \sum_{u_i \in U_I} Y_i$ with Y_i the total of y_k on u_i .

A sample of households is selected by two-stage sampling :

- ullet Sample S_I of municipalities with inclusion probabilities π_{Ii} ,
- Inside any $u_i \in S_I$, subsample S_i of households.

To simplify the presentation, no stratification at the first stage.

Horvitz-Thompson estimator

The Horvitz-Thompson (HT) estimator for the total Y is

$$\hat{Y}_{\pi} = \sum_{u \in S_{-}} \frac{\hat{Y}_{i}}{\pi_{Ii}}$$
 with \hat{Y}_{i} the estimator in the PSU u_{i} .

For first-stage sampling designs with very large entropy, \hat{Y}_{π} is consistent + HT variance estimator consistent (Chauvet and Vallée, 2018).

Simple variance estimators: Hajek-Rosen (HR) or with-replacement (WR)

$$\hat{V}_{HR} = \frac{n_I}{n_I - 1} \sum_{u_i \in S_I} (1 - \pi_{Ii}) \left(\frac{\hat{Y}_i}{\pi_{Ii}} - \frac{\sum_{u_j \in S_I} (1 - \pi_{Ij}) \frac{\hat{Y}_j}{\pi_{Ij}}}{\sum_{u_j \in S_I} (1 - \pi_{Ij})} \right)^2,$$

$$\hat{V}_{WR} = \frac{n_I}{n_I - 1} \sum_{u_i \in S_I} \left(\frac{\hat{Y}_i}{\pi_{Ii}} - \frac{\sum_{u_j \in S_I} \frac{\hat{Y}_j}{\pi_{Ij}}}{n_I} \right)^2.$$

HR: underestimates the variance. WR: overestimates the variance.

Advantage of these simplified variance estimators

They require only total estimators \hat{Y}_i inside PSUs :

- ullet no need to produce unbiased variance estimators \hat{V}_i inside PSUs,
- can be applied with any type of second-stage sampling designs (e.g., systematic sampling),
- ullet consistent when $n_I/N_I
 ightarrow 0$ (Chauvet and Vallée, 2018).

They can be applied if we may (approximatly) write the estimator as a sum of independent random variables :

$$\hat{Y}_{\pi} = \sum_{u_i \in S_I} \frac{Y_i}{\pi_{Ii}}.$$

Important to catch the treatment of unit non-response.

The second variance estimator may be easily bootstrapped.

Bootstrap variance estimation

Computation of bootstrap weights

In case of the Household Finance and Consumption Survey (HFCS), a bootstrap variance estimator is needed.

Bootstrap = computational technique which enables to produce variance estimations by repeatedly mimicking the sampling process + estimation steps.

Proposed bootstrap = bootstrap of Primary Sampling Units (PSUs) :

- does not need to bootstrap the sample selection inside the PSUs
 ⇒ fairly simple to implement,
- ullet overestimation of the 1st-stage sampling variance \Rightarrow conservative,
- variance of further stages (e.g., non-response) correctly accounted for.

Computation of bootstrap weights

The bootstrap weights are obtained as follows :

• Selection of a with-replacement sample of n_I-1 PSUs in the original sample, with equal probabilities (Bootstrap of PSUs).

The bootstrap weight of the PSU u_i is :

$$W_{Ii} = \frac{n_I}{n_I - 1} \times \text{Number of selections of } u_i.$$

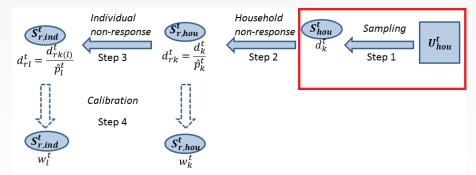
 $oldsymbol{2}$ The bootstrap sampling weight of some household $k \in u_i$ is :

$$d_{k*} = d_k \times W_{Ii}.$$

- These bootstrap sampling weights are adjusted for all the following estimation steps:
 - Household non-response,
 - Individual non-response,
 - Calibration.
- Steps 1-3 are repeated a large number of times (say B=500)

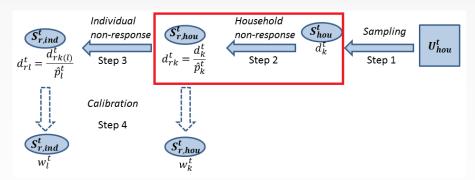


Cross-sectional Estimation Initial time t



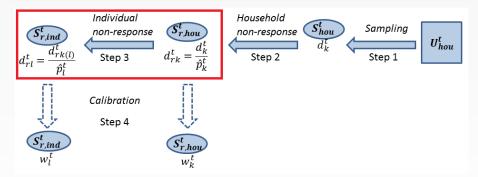
Sample of households selected by multistage sampling : sample of n_I PSUs (e.g., municipalities) and subsampling of households within.

We obtain a sample of households $S_{hou}^{ar{t}}$ with sampling weights d_k^t



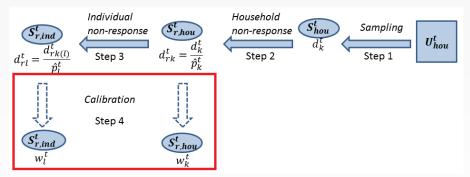
We account for the non-response of households, e.g. via Response Homogeneity Groups (RHGs).

Estimation of the response probability $\hat{p}_k^t \Rightarrow$ household weights d_{rk}^t .



We account for the non-response of individuals, e.g. via Response Homogeneity Groups (RHGs).

Estimation of the response probability $\hat{p}_l^t \Rightarrow$ individual weights d_{rl}^t .



The weights adjusted for non-response are calibrated on auxiliary information on households and individuals.

Leads to calibrated weights w_k^t for households and w_l^t for individuals.

An example with two-stage sampling

Selection of the households

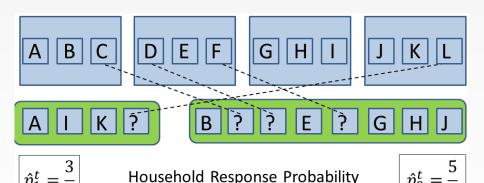
$$\begin{bmatrix} u_1 \\ \pi_{I1} = 0.2 \end{bmatrix} \begin{bmatrix} u_2 \\ \pi_{I2} = 0.4 \end{bmatrix} \begin{bmatrix} u_3 \\ \pi_{I3} = 0.8 \end{bmatrix} \begin{bmatrix} u_4 \\ \pi_{I4} = 0.8 \end{bmatrix}$$

$$\begin{bmatrix} A & B & C \\ D & E & F \\ \end{bmatrix} \begin{bmatrix} G & H & I \\ \end{bmatrix} \begin{bmatrix} J & K & L \\ \end{bmatrix}$$

$$\begin{bmatrix} Household Sampling Weight: d_k^t = \frac{100}{12} \end{bmatrix}$$

Sample of $n_I=4$ PSUs selected by a self-weighted two-stage sampling. We obtain equal sampling weights d_k for the households.

Correction of household non-response



Calibration

$$u_1 \\ \pi_{I1} = 0.2$$

$$\pi_{I2} = 0.4$$

$$u_3 = 0.8$$

$$u_4 \\ \pi_{I4} = 0.8$$

Household Adjusted weights
$$d_{rk}^t = \frac{d_k^t}{\hat{p}_k^t} \Rightarrow$$
 Calibrated weights w_k^t

The sampling weights are corrected to account for non-response.

The weights are finally calibrated.

Variance estimation : bootstrapping of the PSUs

$$u_1 \\ \pi_{I2} = 0.2$$

$$W_{I1} = \frac{4}{3}$$



$$u_4 \\ \pi_{I4} = 0.8$$

$$W_{I4} = \frac{8}{3}$$

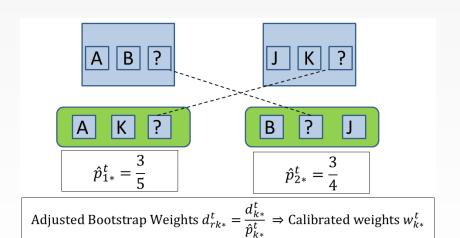
J K ?

Bootstrap Household Sampling Weights $d_{k^*}^t = W_{Ii} \times d_k^t$

A resample of $n_I - 1 = 3$ PSUs is selected with equal probabilities.

The resampling weight of a PSU is 4/3 times its number of selections.

Variance estimation: bootstrapping of non-response and calibration



The bootstrap sampling weights are corrected for non-response as originally.

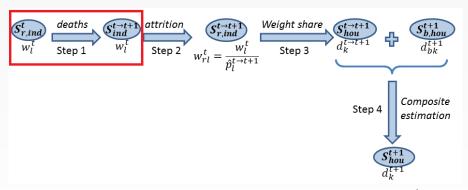
The weights are finally calibrated as originally.

Cross-sectional estimation at time t+1

Panel and refreshment sample

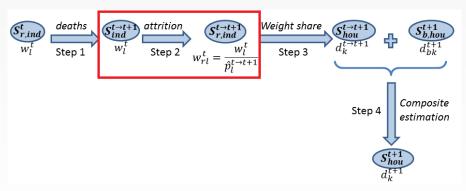


Tracking individuals: elimination of the leaving dead



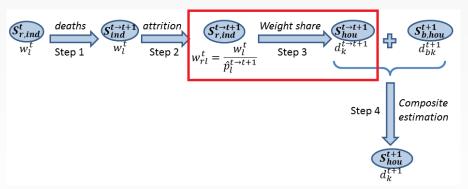
We suppose that individuals are followed over time, with weights $w_l^t.$ We account for deaths between time t and t+1. The weights remain unchanged

Tracking individuals: correction of attrition



We account for attrition, e.g. via Response Homogeneity Groups (RHGs). Estimation of the response probability $\hat{p}_k^{t \to t+1} \Rightarrow$ individual weights w_{rl}^t .

From individual to household: Weight Share Method

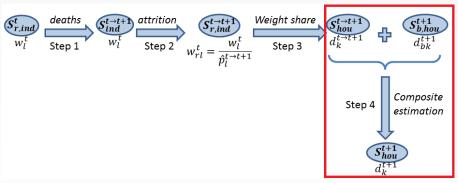


The individual weights are distributed on their household via the Weight Share Method.

This leads to the household weights $d_k^{t o t+1}$.



Merging the panel and the refreshment sample

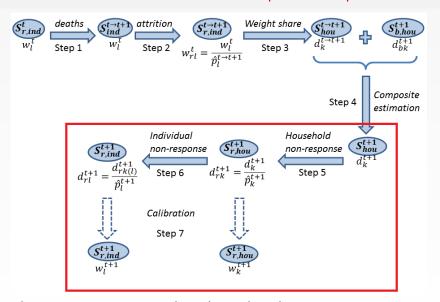


The panel sample and the refreshment sample are joined, by using composite estimation.

This leads to the household weights d_k^{t+1} .

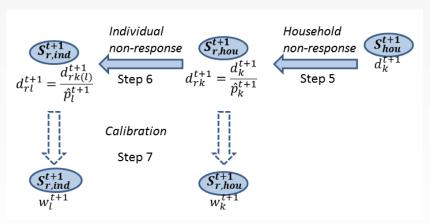


Cross-sectional Estimation with a panel component



The estimations steps are then identical to that at previous time t. Accounting for non-response of households + individuals and calibration.

Cross-sectional Estimation with a panel component



Bootstrap weights computed as if S_{hou}^{t+1} was directly selected at time t+1. Expected to overestimate slightly the variance due to attrition and the variance due to household non-response at time t. Other variance components correctly accounted for.

4D > 4A > 4B > B 990

Related/further work

Related/further work

Longitudinal estimation can be handled similarly (slightly easier).

Need for a software: macro SAS in progress. Interested for applications to other surveys.

Theoretical justification: works (fairly) well in simulations, but the properties of inverse probability weighted estimators used in case of unit non-response are more complicated to handle.



A small simulation study

Simulated data

Original sampling design of the PPV survey :

- stratified sample of districts (PSUs),
- inside districts, sample of households (SSUs).

Responding households and districts duplicated by inverse sampling rates \Rightarrow artificial but realistic population :

- ullet 433 districts, ranging from 600 to 23 000 households,
- population of 1.1 Mo households.

Variables of interest:

- reputation of district (good or very good/other),
- witness of trafficking (never or rarely/other),
- urban works during last year (yes/no),
- move planned during year to come (yes/no).



Simulated sampling design

Unstratified self-weighted two-stage design :

- Inclusion probabilities π_{Ii} proportional to size $(n_I = 20 \text{ or } n_I = 40)$, \Rightarrow with $n_I = 20$, π_{Ii} ranging from 0.01 to 0.42 (mean 0.05).
- ullet Sample of $n_0=50$ households inside each selected district.
- 6 Response Homogeneity groups obtained by crossing :
 - low-income housing (yes/no),
 - region (Ile-de-France/north/south).

Average response probability : 74%.

Sampling design + non-response repeated $B=5\ 000$ times to obtain a Monte Carlo variance approximation.

Simplified variance estimators computed ${\cal C}=500$ times.

Relative bias of the simplified variance estimators

	\hat{V}_{SIMP1}		\hat{V}_{SIMP2}	
	$n_I = 20$	$n_I = 40$	$n_I = 20$	$n_I = 40$
reputation	-5.7 %	-3.2 %	+3.8 %	+18.2 %
trafficking	-4.3 %	-0.9 %	+6.3 %	+23.5 %
works	+0.4 %	-3.6 %	+7.2 %	+10.4 %
move	-3.8 %	-9.8 %	+6.5 %	+12.3 %

Relative bias of the simplified variance estimators

	\hat{V}_{SIMP1}		\hat{V}_{SIMP2}	
	$n_I = 20$	$n_I = 40$	$n_I = 20$	$n_I = 40$
reputation	-5.7 %	-3.2 %	+3.8 %	+18.2 %
trafficking	-4.3 %	-0.9 %	+6.3 %	+23.5 %
works	+0.4 %	-3.6 %	+7.2 %	+10.4 %
move	-3.8 %	-9.8 %	+6.5 %	+12.3 %