c⸈ C. 2

Methodology Branch

Household Survey
Methods Division

Direction de la méthodologie

Division des méthodes
d'enquêtes des ménages

Canadä

# CROSS-SECTIONAL ESTIMATION IN MULTIPLE-PANEL HOUSEHOLD SURVEYS

HSMD - 99 - 004E

Takis Merkouris

Household Survey Methods Division
Statistics Canada

October 1999

# Cross-sectional Estimation in Multiple-Panel Household Surveys

TAKIS MERKOURIS [1]

## ABSTRACT

This paper presents weighting procedures that combine information from multiple panels of a repeated panel household survey for cross-sectional estimation. The dynamic character of a repeated panel survey is discussed in relation to estimation of population parameters at any wave of the survey. A repeated panel survey with overlapping panels is described as a special type of multiple frame survey, with the frames of the panels forming a time sequence. The paper proposes weighting strategies suitable for various multiple panel survey situations. The proposed weighting schemes involve an adjustment of weights in domains of the combined panel sample that represent identical time periods covered by the individual panels. A weight adjustment procedure that deals with changes in the panels over time is discussed. The integration of the various weight adjustments required for cross-sectional estimation in a repeated panel household survey is also discussed.

KEY WORDS: Repeated panel surveys; Multiple frames; Temporal domains; Combined panels; Cross-sectional weighting; Weight share method.

[1] Takis Merkouris, Household Survey Methods Division, Statistics Canada, Ottawa, Ontario, K1A 0T6.

# L'estimation transversale dans les enquêtes-ménages à panels multiples

TAKIS MERKOURIS [2]

## RÉSUMÉ

Ce document décrit des procédures de pondération qui combinent les renseignements provenant de panels multiples d'enquêtes-ménages à passages répétés en vue de produire une estimation transversale. Le caractère dynamique des enquêtes par panels à passages répétés est étudié par rapport à l'estimation des paramètres de population à n'importe quel cycle des enquêtes. L'enquête à panels chevauchants et à passages répétés est décrite comme un genre particulier d'enquête à bases de sondage multiples, ces dernières formant une suite chronologique. Le document propose des stratégies de pondération qui conviennent à diverses formes d'enquêtes à panels multiples. Les procédés de pondération proposés exigent un ajustement des poids dans les domaines de l'échantillon de panels combinés qui représentent des périodes identiques couvertes par des panels distincts. Le document décrit une procédure d'ajustement des poids qui tient compte des changements survenus au sein des panels avec le temps. Il explique aussi l'intégration des divers ajustements de poids nécessaires à l'estimation transversale dans des enquêtes-ménages par panels à passages répétés.

MOTS CLÉS : enquêtes par panels à passages répétés, bases de sondage multiples, domaines temporels, panels combinés, pondération transversale, méthode de partage de poids.

[2] Takis Merkouris, Division des méthodes d'enquêtes auprès des ménages, Statistique Canada, Ottawa, Ontario, K1A 0T6.

# 1. INTRODUCTION

A panel survey collects the survey data for the same sample elements at different time points (the survey waves). A repeated panel survey is made up of a series of panel surveys, each having fixed duration, with the panels selected at different time points. In a repeated panel household survey a sample of households is selected for each panel from the population of households existing at the start of the panel. All the individuals in the sampled households become panel members to be followed throughout the duration of the panel or until they leave the survey population. At a subsequent survey wave the household sample consists of all the households in which panel members reside. A review of various types of panel surveys is given in Kalton and Citro (1993). A formalization of related concepts can be found in Deville (1998).

The type of repeated panel household survey considered in this paper consists of overlapping panels, with two or more panels covering overlapping fractions of the same time period. A typical example of such a survey is the Canadian Survey of Labour and Income Dynamics (SLID), which employs two overlapping panels of duration of six years each; for a description of SLID see Lavigne and Michaud (1998). In SLID, each new panel is introduced three years after the introduction of the previous one. The sample for each panel is made up of two rotation groups from the Canadian Labour Force Survey, which uses a stratified multistage design with an area frame wherein dwellings containing households are the final sampling units.

A panel survey, though primarily conducted for longitudinal purposes, may also be used to produce cross-sectional estimates of population parameters for any survey wave. For cross-sectional purposes, data are usually collected at each survey wave for all individuals living in households that contain at least one selected member. The process of obtaining cross-sectional estimates at any wave of a panel household survey after the first wave presents difficulties arising from the dynamic character of the panels. Weighting schemes that deal with dynamic features of a single panel, such as movers and "cohabitants," have been discussed in the literature; see Kalton and Brick (1995), and Lavallée (1995) for details. Yet, there seems to be a paucity of work in the literature on cross-sectional weighting and estimation for repeated panel household surveys with overlapping panels. Some initial work in the context of SLID can be found in Lavallée (1994).

This paper describes procedures for cross-sectional estimation that combine information from overlapping panels of a repeated panel household survey. The coverage of the population at any given wave by the individual panels, and the use of the combined panels supplemented by a "top-up" sample to construct a representative cross-sectional sample are discussed in Section 2. Also discussed in the same section are the analogies with a multiple-frame survey scheme, as well as issues related to the dynamic character of the sample. The weighting and estimation problem in repeated panel household surveys is described in Section 3. Weighting strategies suitable for various panel survey situations are then proposed. Bias and efficiency issues related to the combination of panels are discussed. A weight adjustment procedure that deals effectively with changes in the combined panels over time is described in Section 4. The integration of the various weight adjustments required in cross-sectional estimation for a repeated panel household survey is discussed in Section 5. Concluding remarks on the proposed procedures are made in Section 6.

# 2. GENERAL CONSIDERATIONS

## 2.1 Coverage of the cross-sectional population

Important to cross-sectional estimation are changes in the population composition over time, occurring when individuals leave or enter the population. In a single-panel household survey, new entrants who have joined the survey population since the start of the panel are not represented in the sample at later waves if they live in households that do not contain any members of the original population. A household survey with

multiple overlapping panels provides a better coverage of the survey population, as it reduces the time period not covered by any of the panels. In the case of SLID, this time period is reduced from a maximum of six years to a maximum of three years. Nevertheless, the problem of complete coverage remains, unless a special supplementary sample of the non-covered population is taken at each survey wave. A survey scheme involving one panel and a supplementary sample drawn at each survey wave for cross-sectional purposes is described in Lavallée (1995). An alternative approach involves the selection, at each wave, of a new sample that covers the entire survey population but does not form a new panel. This sample (henceforth to be called top-up) is to be used only once, for cross-sectional purposes, and its size would normally be smaller than a panel's size. Thus, a household survey with multiple overlapping panels and a top-up sample at each wave provides complete coverage of the target population for cross-sectional purposes.

The situation with regard to individuals who leave the population is as follows. For any panel, the sampling frame for the survey population at a time point $t$ is essentially the sampling frame for the population at the start of the panel, with the leavers in the intervening period being treated as blanks on the frame. Panel members who leave the population before time $t$ correspond to blanks on the frame, and thus their effect on cross-sectional estimates at time $t$ is loss of efficiency but not bias; see also Kalton and Brick (1995) for relevant discussion.

The foregoing observations lead to the following perspective regarding the coverage of the population by each of the panels at any wave of the survey. As regards cross-sectional representation, each panel covers at the time of its selection the entire survey population represented by the preceding panels. Accordingly, the frames of the panels form a time sequence, with the frame of each panel containing at the start of the panel the frames of the preceding panels. In such a sequence of frames a common (overlap) frame is formed sequentially as the intersection of the frame of a new panel with the remaining of the original common frame of the preceding active panels. At any wave the common panel frame is the common frame at the start of the most recent of these panels, but without the leavers. The non-overlap frame domain at the start of a new panel consists of individuals who entered into the population after the start of the preceding panel. Other frame domains (relatively very small in size) may be formed by returning units of older frames, in which case the time sequence of frames is not completely nested. Because of the latter type of frame domains the complete frame at any wave after the selection of the most recent panel is the union of the frames of all panels at that time point, no just the remaining of the frame of the most recent panel. In panel surveys that employ a top-up sample at each wave, essentially as a small panel, the complete frame is that of the top-up sample.

## 2.2 A multiple frame analogy

With the above considerations, a multiple panel survey with overlapping panels can be thought of as a special type of multiple frame survey, in which the frame for the cross-sectional population is the union of mutually exclusive temporal domains defined by the frames of the panels and their intersections. The sizes of the frames of the individual panels as well as the characteristics of the population members in each panel's frame change over time. This is in contrast with the static character of the usual type of multiple frame survey. Also, there is a high degree of nesting in the sequence of panel frames, so that the total number of mutually exclusive temporal frame domains is small. Among the various frame domains the one that is common to all panels is by far the largest. These special multiple frame features have implications in cross-sectional estimation, as will be discussed in the next section.

The sample temporal domains may be more dynamic because of attrition, moves of selected individuals within and between panels and moves of non-selected individuals into households in which panel members reside. For instance, with the presence of new entrants (e.g., immigrants) in households that contain selected individuals, a panel crosses the boundary of its frame into the frame of the succeeding panel.

-2-

The analogy with multiple-frame survey sampling places the problem of cross-sectional estimation for repeated surveys with overlapping panels into a familiar framework. However, the distinctive dynamic features of multiple panel surveys will have to be considered if conventional multiple frame approaches are contemplated for the formulation of a cross-sectional estimation methodology. The difficulty in developing a cross-sectional estimation procedure for multiple panel surveys with overlapping panels arises from the complications that the changes in the population and in the sample add to the more standard problem of combining information from multiple sources.

For the purpose of introducing a cross-sectional estimation procedure that combines information from the panels of a repeated panel household survey, it suffices to consider the simple situation involving two overlapping panels at the time point of the start of the second panel. Note that this would always be the situation in a survey with one panel and a top-up sample. Thus, using multiple frame notation, with $B$ and $A$ denoting the frames of the first and the second panel ($B \subset A$) at the start of the second panel, and with $s_B$, $s_A$ denoting the respective samples, the setting can be presented schematically as in Figure 1.
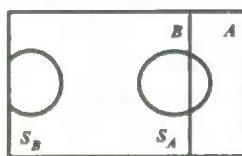


Figure 1. Two overlapping panels at
the start of the second panel.

In Figure 1, $A$ is the complete frame, so that at its start the second panel represents the cross-sectional population at that time. The overlap domain $B$ is the remaining of the original frame of the first panel. The domain $a = B^c \cap A$ consists of all new entrants into the population since the start of the first panel. The samples $s_B$ and $s_A$ are the originally selected ones, with $s_B$ reduced in size because of leavers and non-respondents. It is assumed that the samples $s_A$ and $s_B$ are drawn independently from $A$ and $B$ according to specified probability designs $p_A(s_A)$ and $p_B(s_B)$, which determine the inclusion probabilities $\pi_{Ai}$ and $\pi_{Bi}$ of the $i$-th unit (household or any individual within it) for the original samples $s_A$ and $s_B$, respectively. The samples $s_A$ and $s_B$ may intersect, since members in the overlap frame $B$ can be selected in both panels. The issue of panel (sample) overlap is akin to that of duplicate sample units in multiple frame surveys. In repeated panel household surveys an operational constraint motivated by respondent burden may be to exclude from $s_A$ individuals already selected in $s_B$, thus inducing $s_A \cap s_B = \varnothing$; for a discussion on this see Lavallée (1994). Here, as in any multiple frame situation, it is observed that if the probabilities $\pi_{Ai}$ and $\pi_{Bi}$ are small the probability of duplicate units is negligible. It will be assumed in the following that the probabilities $\pi_{Ai}$ and $\pi_{Bi}$ are small, and in effect $s_A \cap s_B = \varnothing$.

## 3. CROSS-SECTIONAL WEIGHTING AND ESTIMATION

This section describes procedures that combine information from multiple panels of a repeated panel household survey for cross-sectional estimation of population parameters. The discussion is confined to estimation of totals. A uniform approach to cross-sectional estimation for households and individuals is presented. This approach is based on the production of a set of weights for the combined panel sample that yield design-unbiased estimators of cross-sectional totals. Essentially, it involves the construction of a combined cross-sectional sample by means of an adjustment of the sampling weights of units from the

temporal domains of the different panels that represent common temporal domains of the cross-sectional population. While the identification of the various temporal frame domains is necessary for determining the coverage of parts of the cross-sectional population by different panels, the identification of some of the corresponding sample domains may not be possible under the operating procedures of a repeated panel household survey. For example, the information needed to determine whether or not a unit in the second panel belongs to the non-overlap frame domain $a$ (see Figure 1) may not be available. In this section, both cases of identifiable and non-identifiable temporal sample domains are considered. A "weight share" adjustment that handles changes in the sample composition over time is to follow the combination of the panels, as it can be applied readily only to the combined sample; see relevant discussion in Section 4.

### 3.1 Identifiable temporal sample domains

For the construction of a cross-sectionally representative combined sample, a panel survey scheme as that depicted in Figure 1 is considered. The two samples $s_A$ and $s_B$ can be thought of as selected independently from the complete frame $A$, but with a fixed time lag between the two selections, according to the sampling designs $p_A(s_A)$ and $p_B(s_B)$. The two sampling designs $p_A(s_A)$ and $p_B(s_B)$ induce a well-defined design $p(s)$ on the set of samples $s = s_A \cup s_B$ in $A$. Thus conventional estimators, based on a single frame and a combined sample, may be constructed from $p(s)$. The standard approach, leading to the Horvitz-Thompson estimator, would be to assign sample units weights made inversely proportional to their inclusion probabilities. The inclusion probability $\pi_i = P(i \in s)$ of the $i$-th unit of the combined sample $s$ is $\pi_{Ai} + \pi_{Bi} - \pi_{Ai}\pi_{Bi}$ if $i \in s \cap B$, and $\pi_{Ai}$ if $i \in s \cap a$. The weight of the $i$-th unit of the sample is then $w_i = 1/\pi_i$. This weighting scheme can be used provided that it is possible to identify the common units in the samples $s_A$ and $s_B$, so that the duplicate units can be eliminated. A simpler approach, especially for surveys with more than two panels, would be to assign sample units in $s \cap B$ weights made inversely proportional to their expected number of selections, that is, inversely proportional to $\pi_{Ai} + \pi_{Bi}$. This weighting scheme, proposed by Kalton and Anderson (1986) for multiple frame surveys, does not require identification of duplicate sample units. Now, consider the domains $s_{ab} = s_A \cap B$ and $s_a = s_A \cap a$ of $s_A$. Also, let a value $y_i$ be associated with population unit $i$ for some population characteristic, and define the population total $Y_A = \sum_A y_i \; (= \sum_B y_i + \sum_a y_i)$. Then, employing the latter weighting scheme, the unbiased estimator

$$\hat{Y}_A = \sum_s w_i y_i = \sum_{s_B \cup s_{ab}} (\pi_{Ai} + \pi_{Bi})^{-1} y_i + \sum_{s_a} \pi_{Ai}^{-1} y_i \tag{1}$$

of the total $Y_A$ can be constructed. On the assumption that the probabilities $\pi_{Ai}$ and $\pi_{Bi}$ for $i \in s \cap B$ are small, the estimator $\hat{Y}_A$ is approximately equal to the Horvitz-Thompson estimator.

The approach leading to the estimator (1) is not in general feasible, since determination of the weight $w_i = \pi_{Ai} + \pi_{Bi}$ for $i \in s \cap B$ requires that the inclusion probabilities of the sampled units be known over both frames, which is difficult or impossible to ascertain in household surveys. In multiple-panel household surveys additional complications arise from the time element. For units that move (e.g., to another stratum) in the time between the selection of the panels it is impossible to determine both $\pi_{Ai}$ and $\pi_{Bi}$.

An alternative strategy needs to be considered for developing weights for the sample overlap domain $s \cap B$. An approach that provides a general framework for handling this problem requires information on the probability of inclusion in only one of $s_A$ or $s_B$, thus avoiding the difficulty noted above. The essence of the alternative approach considered here is to associate with the $i$-th unit from the overlap frame $B$ a number $p_i$ ($0 \le p_i \le 1$) when the unit is selected in $s_B$, and the number $1 - p_i$ when the unit is selected in $s_A$, and then define the weight of the unit as

$$w_i^* = p_i \frac{1}{\pi_{Bi}} I\{i \in s_B\} + (1-p_i)\frac{1}{\pi_{Ai}} I\{i \in s_{ab}\} \ , \quad i \in B \ , \tag{2}$$

where $I$ is the usual sample membership indicator variable. Clearly, $E(w_i^*) = 1$ under $p(s)$, and thus the use of the weights $w_i^*$ will yield unbiased estimators $\hat{Y}_B = \Sigma_B w_i^* y_i$ for the total $Y_B = \Sigma_B y_i$, for any choice of constants $p_i$ satisfying $0 \le p_i \le 1$, and for any sampling designs $p_A(s_A)$ and $p_B(s_B)$. Equation (2) can be written alternatively as $w_i^* = p_i w_{Bi} + (1-p_i) w_{Ai}$, with the obvious definition of the weights $w_{Bi}$ and $w_{Ai}$ associated with the samples $s_B$ and $s_A$. Thus, the class of weighting schemes defined by equation (2) consists essentially of different weighted combinations of the weights in the original samples $s_B$ and $s_A$. The limits on the values of $p_i$ ensure that the weight $w_i^*$ will be nonnegative. Note that the intractable weight $w_i = (\pi_{Ai} + \pi_{Bi})^{-1}$, for $i \in s \cap B$, used in (1) is a special case of $w_i^*$ with $p_i = \pi_{Bi}(\pi_{Ai} + \pi_{Bi})^{-1}$.

Evidently, the weighting scheme defined by (2) does not eliminate duplicate units that fall in both samples. If the operational constraint to exclude from $s_A$ individuals already selected in $s_B$ is imposed, the second term in the right-hand side of (2) should be modified to $(1-p_i)[\pi_{Ai}(1-\pi_{Bi})]^{-1} I\{i \in s_{ab}, i \notin s_B\}$ to ensure that $E(w_i^*) = 1$. This, however, may be impossible to do since it requires that the inclusion probabilities of the sampled units be known over both frames. Note also that under the constraint of excluding duplicate units the two samples will not be independent. Nevertheless, as it is assumed that both probabilities $\pi_{Ai}$ and $\pi_{Bi}$ are small, the probability of duplicate units will be negligibly small, and hence any bias effect resulting from using the tractable weighting scheme defined by (2) would also be negligible. On this assumption, the two indicator variables in (2) should be understood to satisfy $I\{i \in s_B\} I\{i \in s_{ab}\} = 0$.

The question arises now as to an optimal choice of $p_i$, for any $i \in s \cap B$, according to some criterion of optimal weighting for the combined sample. One approach is to choose the $p_i$ to minimize the variance of the estimated total $\hat{Y}_A = \Sigma_B w_i^* y_i + \Sigma_a w_i y_i$, where $w_i = (\pi_{Ai})^{-1} I\{i \in s_a\}$ for $i \in a$. However, minimization of the variance of $\hat{Y}_A$ with respect to $p_i$ for all $i \in s \cap B$ is not tractable. A simpler option is to restrict the class of weighting schemes defined by equation (2) to one in which the weight adjustment factors are specified not at the unit level but rather at a higher level, which may be a stratum or the entire overlap frame $B$. Further discussion on the level of adjustment is deferred until the end of this subsection. It suffices for the development of the weighting procedure to consider next the case involving a uniform weight adjustment factor $p$ for the entire frame $B$. Then, the class of weighting schemes defined by equation (2) for the frame domain $B$ generates a class of unbiased estimators for the overall total $Y_A$ of the form

$$\hat{Y}_A^p = p\hat{Y}_{s_B} + (1-p)\hat{Y}_{s_{ab}} + \hat{Y}_{s_a}, \tag{3}$$

where $\hat{Y}_{s_B}$ and $\hat{Y}_{s_{ab}}$ are independent Horvitz-Thompson estimators of $Y_B$ based on $s_B$ and $s_{ab}$, respectively, and $\hat{Y}_{s_a}$ is the Horvitz-Thompson estimator of $Y_a$ based on $s_a$. The limits on $p$ ensure that $\hat{Y}_A^p$ will be nonnegative whenever the $y_i$ are nonnegative. The limit values of $p$ yield two special cases of the estimator $\hat{Y}_A^p$, in both of which the overlap domain total $Y_B$ is estimated from one panel only. When $p$ is set equal to zero in (3), the resultant trivial estimator $\hat{Y}_A^p$ for the entire population is based only on $s_A$. More notable is the case with $p$ set equal to one in (3). The implied simple unbiased estimator $\hat{Y}_A = \hat{Y}_{s_B} + \hat{Y}_{s_a}$ would be the natural estimator in a panel survey with one panel and a supplementary sample from the population of new entrants, with the units in the supplementary sample being "screened," and only the units in the domain of new entrants being enumerated. In such a context this simple estimator would be a special case of a "screening" multiple frame estimator, the special feature being the temporal nature of the non-overlap frame domain $a$. In the present context the screening estimator appears inefficient because information in the sample $s_{ab}$ is not utilized. Better use may be made of data from both panels by combining $s_B$ and $s_{ab}$, using

-5-

an optimal $p$ that is based on the minimization of the variance of $\hat{Y}_A^p$. The optimal value of $p$ is given by

$$p = \frac{Var(\hat{Y}_{s_{ab}}) + Cov(\hat{Y}_{s_{ab}}, \hat{Y}_{s_a})}{Var(\hat{Y}_{s_B}) + Var(\hat{Y}_{s_{ab}})}. \tag{4}$$

The variance and covariance terms in (4) are unknown, but could be estimated from the sample data, in which case the chosen $p$ would actually minimize the estimated variance of $\hat{Y}_A^p$. There are numerous problems associated with this choice of value for $p$, the obvious one being that estimation of the optimal $p$ is inconvenient, especially in surveys involving more than two panels. Furthermore, the dependency of the estimated optimal $p$ on the sample data entails $E(w_i^*) \neq 1$ for $i \in B$, which disturbs the unbiasedness of the estimator (3). It is to be noted that the condition $E(w_i^*) = 1$ is also necessary for the validity of the weight share method (see Section 4) to hold when applied to the combined sample $s$ at any wave after the selection of the second panel. It is also noted here, in passing, that a sample estimate of the optimal $p$ in (4) would add variability into the estimator $\hat{Y}_A^p$, and complicate the estimation of its variance.

An alternative choice for the value of $p$ is based on the minimization of the variance of the estimator of the common-frame total $\hat{Y}_B^p = p\hat{Y}_{s_B} + (1-p)\hat{Y}_{s_{ab}}$. This restricted minimization, which ignores the typically small domain estimator $\hat{Y}_{s_a}$, gives the value

$$p' = \frac{Var(\hat{Y}_{s_{ab}})}{Var(\hat{Y}_{s_B}) + Var(\hat{Y}_{s_{ab}})}, \tag{5}$$

which is independent of the covariance term, and always lies between zero and one. If the variance of $\hat{Y}_B^p$ conditional on the realized value of the random size $n_{ab}$ of the sample domain $s_{ab}$ is minimized, and if finite population corrections are disregarded, then it can be shown that (5) may be written as

$$\hat{p}' = \frac{n_B d_{ab}}{n_B d_{ab} + n_{ab} d_B}, \tag{6}$$

where $n_B$ is the size of the sample $s_B$, and $d_B$, $d_{ab}$ are the design effects associated with the samples $s_B$ and $s_{ab}$. The calculation of the value of $\hat{p}'$ requires estimates of the two design effects, which need not be based on the samples $s_B$ and $s_{ab}$. Suitable approximate values of $d_B$ and $d_{ab}$ may be available from other surveys with the same sampling designs as the two panels. However, the dependency of $\hat{p}'$ on the variable $y$ through $d_B$ and $d_{ab}$ requires a compromise solution. To this end, the approximate values of $d_B$ and $d_{ab}$ could preferably be obtained for a count variable associated with a large portion of the survey population and correlated with main survey variables. It is to be noted that since $\hat{p}'$ depends on the characteristic $y$ only through the ratio $d_B/d_{ab}$, the loss of efficiency for estimators of totals of other characteristics should not be substantial. It is to be noted further that because of the time lag between the selection of the two panels, the design effects will be different, and thus present in (6), even when the sampling designs for the two panels are identical. By using estimates of the design effects from external sources the randomness of $\hat{p}'$ is only due to the random size of the sample domain $s_{ab}$. Since the size of the sample $s_A$ is usually very large, and the size of the overlap frame $B$ is typically only a little smaller than the size of the complete frame $A$, the size $n_{ab}$ of the sample domain $s_{ab}$ must be nearly constant, and thus the unbiasedness condition $E(w_i^*) = 1$ will hold approximately.

Some loss of efficiency will be incurred by ignoring $\hat{Y}_{s_a}$ in deriving an optimal value for $p$, but this loss may be insignificant given the relatively very small size of the domain $a$ in most household panel surveys, because of the typically small time lag between panels. To assess this loss of efficiency, let $\hat{Y}_A$ and $\hat{Y}'_A$ denote the estimator $\hat{Y}^p_A$ when the values of $p$ given by expressions (4) and (5), respectively, are substituted in expression (3). Then, a simple calculation gives

$$Var(\hat{Y}'_A) - Var(\hat{Y}_A) = \frac{Cov^2(\hat{Y}_{s_{ab}},\hat{Y}_{s_a})}{Var(\hat{Y}_{s_B}) + Var(\hat{Y}_{s_{ab}})} \leq \frac{Var(\hat{Y}_{s_{ab}})Var(\hat{Y}_{s_a})}{Var(\hat{Y}_{s_B}) + Var(\hat{Y}_{s_{ab}})} = p'Var(\hat{Y}_{s_a}),$$

so that an upper bound for the efficiency loss can be obtained as

$$\frac{Var(\hat{Y}'_A) - Var(\hat{Y}_A)}{Var(\hat{Y}_A)} \leq p'\frac{Var(\hat{Y}_{s_a})}{Var(\hat{Y}_A)}.$$

Given the usually very small size of $\hat{Y}_{s_a}$ relative to $\hat{Y}_A$ (the size of the domain $a$ may be as small as one fortieth of the size of the complete frame $A$ in the case of SLID) it appears that the loss of efficiency will be very small in most panel household surveys.

An interesting question is whether or not $\hat{Y}_A$ is more efficient than the simple "screening" estimator $\hat{Y}_A = \hat{Y}_{s_B} + \hat{Y}_{s_a}$, whose variance is $Var(\hat{Y}_{s_B}) + Var(\hat{Y}_{s_a})$. It can be readily shown that $Var(\hat{Y}'_A) < Var(\hat{Y}_{s_B}) + Var(\hat{Y}_{s_a})$ if $2Cov(\hat{Y}_{s_{ab}},\hat{Y}_{s_a}) < Var(\hat{Y}_{s_B})$, which certainly holds if the covariance of $\hat{Y}_{s_{ab}}$ and $\hat{Y}_{s_a}$ is negative. This covariance may in fact be positive, given that $s_{ab}$ and $s_a$ are domains of the sample $s_A$ which is drawn from an area frame according to a cluster design. In that case too, however, the condition will most likely hold, given the magnitude of $Var(\hat{Y}_{s_B})$ relative to $Var(\hat{Y}_{s_a})$, and the magnitude of $Var(\hat{Y}_{s_{ab}})$ relative to $Var(\hat{Y}_{s_a})$. Indeed, the sizes of the panel samples $s_B$ and $s_A$ are typically equal by design, although the effective panel sizes (realized sizes, at any wave, adjusted for design effects) may be considerably different due to differential attrition and design effects between the two panels. Also, with the sizes of the sample domains $s_{ab}$ and $s_a$ roughly proportional to the corresponding population domain sizes, $Var(\hat{Y}_{s_a})$ will be many times, say $k$, smaller than $Var(\hat{Y}_{s_{ab}})$. Then,

$$2Cov(\hat{Y}_{s_{ab}},\hat{Y}_{s_a}) \leq 2\sqrt{Var(\hat{Y}_{s_{ab}})Var(\hat{Y}_{s_a})} = 2\frac{Var(\hat{Y}_{s_{ab}})}{\sqrt{k}},$$

so that at a sufficient condition for the estimator $\hat{Y}'_A$ to be more efficient than the "screening" estimator is

$$2\frac{Var(\hat{Y}_{s_{ab}})}{\sqrt{k}} < Var(\hat{Y}_{s_B}).$$

The interpretation of this is that the sample domain $s_{ab}$ is not to be ignored when estimating $Y_A$ if $Var(\hat{Y}_{s_a})$ is not too small relative to $Var(\hat{Y}_{s_{ab}})$. The condition is ordinarily satisfied in panel household surveys. An additional argument in favour of including $s_{ab}$ in estimation is its better quality relative to $s_B$, since the latter is more liable to the potential bias effect of sample attrition.

The simple approximate weight adjustment factor $\hat{p}'$ given by expression (6) affords an efficient combination of panel samples, accounting for the precision of $\hat{Y}_{s_B}$ relative to that of $\hat{Y}_{s_{ab}}$ through the effective

sample sizes $n_B/d_B$ and $n_{ab}/d_{ab}$. These effective sample sizes are time-dependent, though their ratio (and hence $\hat{p}'$) should be quite stable over the period of panel overlap. Regarding variance calculations, since $n_{ab}$ is nearly non-random, the adjustment factor $\hat{p}'$ can be conveniently treated as constant in any variance estimation procedure.

It is important to emphasize here that additional gains in efficiency will result from the incorporation of auxiliary information into the weights through a calibration weight adjustment to known population totals.

If the criterion in the choice of the value of $p$ is the minimization of the mean square error of the estimator of the common-frame total $\hat{Y}_B^P = p\hat{Y}_{s_B} + (1-p)\hat{Y}_{s_{ab}}$, then it can be easily shown that when the biases of $\hat{Y}_{s_B}$ and $\hat{Y}_{s_{ab}}$ are equal the optimal value of $p$ is the same as the one given by (5). The biases (if any) are not expected though to be equal; for instance, the differential sample attrition rates for the two panels may result in different levels of biases. It is clear that the bias of the linear combination $\hat{Y}_B^P = p\hat{Y}_{s_B} + (1-p)\hat{Y}_{s_{ab}}$, though not minimized if $p$ is as in (5), is nevertheless smaller than the larger of the two component biases. Other complexities aside, the unavailability of good estimates for the two biases renders the criterion of minimum mean square error intractable.

The weighting procedure described above applies to the simple situation of a two-panel survey at the start of the second panel. At later survey waves an additional non-overlap frame domain, denoted by $b$, may be formed by returning leavers of the frame $B$. Units from $b$ originally selected in the first panel were not present when the second panel was selected. Clearly, the weights in the non-overlap sample domain $s_b$ are not to be adjusted for the purpose of combining the two panels. Furthermore, the value for $p$ will not be affected, as it is based only on the overlap part of the combined sample. As with ignoring the sample $s_a$ in determining the value of $p$, ignoring the much smaller, possibly void, sample $s_b$ will have negligible impact on the efficiency of derived estimators.

The simplicity of the proposed weighting procedure for the combination of two panels makes its generalization to situations involving more than two overlapping panels straightforward. The construction of an efficient combined cross-sectional sample would then involve the adjustment of the sampling weights of units from temporal domains of the different panels that represent common temporal domains of the cross-sectional population. For each temporal population domain the weight adjustment factors will be based on the relative effective sample sizes of the corresponding panel domains, in analogy with expression (6), and they will add up to one. The total number of mutually exclusive temporal frame domains is small, because of the high degree of nesting in the sequence of panel frames, and thus the calculations required in the determination of the corresponding independent sets of adjustment factors will not be excessive.

Alternative estimation techniques from the general theory of multiple frame surveys with complex designs (for an account, see Skinner and Rao (1996), and Singh and Wu (1996)) would not be preferable in the present context, for reasons similar to those stated in the discussion ensuing equation (4), or even applicable in surveys with more than two panels. Moreover, when applied to a multiple panel survey such techniques would ordinarily involve the construction of cross-sectional estimators as linear combinations of estimators derived separately from temporal domains of the different panels that represent common temporal domains of the cross-sectional population. Although such cross-sectional estimators would have the same form as those used in this paper for the intermediate purpose of combining the panels, their components from each panel would be poststratified or, more generally, calibrated estimators based on the sample of that panel only. As such, these panel estimators would have incorporated all the weight adjustments, including the "weight share" adjustment, separately for each panel. This would be in conflict with the proposition that the "weight share" adjustment should be applied to the combined sample. It is interesting to note that apart from this complication there are many possible limitations that could have rendered a separate calibration of each panel unfeasible. Membership in the different temporal domains may not be possible to determine for all sample units. Auxiliary totals of frames of old panels that account for the loss of population units may not be available. Many separate calibrations, for each domain of each panel, would be required. Furthermore, since

all temporal sample domains (except the one that is common to all panels) are typically very small, a calibration involving a large number of auxiliary totals (as customary in household surveys) would not be sensible for reasons of potential bias and loss of efficiency of derived estimators of characteristics of interest. It should also be pointed out that accurate auxiliary totals would most likely be unavailable if the frame of each panel were augmented with new entrants who live with individuals of the original frame of the panel. Such would be the situation if the "weight share" procedure, which assigns a basic weight to new entrants living with selected individuals, were to precede the combination of the panels. Lastly, a known drawback of various multiple frame techniques is that a different set of weights would need to be calculated for each characteristic of interest. Besides making the estimation process operationally very inconvenient, the different sets of weights may lead to inconsistencies among estimates.

Returning now to an earlier point, varied weight adjustment factors may be specified at a lower level of sample grouping, such as a certain stratification level. For reasons of feasibility (identical stratification for the two panels is required for that level) and operational convenience, a high level of stratification should be chosen. The natural choice is a superstratum level, at which all other weighting and estimation procedures are carried out independently for each superstratum. Such superstrata could be states or, as in the case of SLID, provinces. The advantage of specifying weight adjustment factors at the province, say, level is improved efficiency, since an optimal or nearly optimal weight adjustment factor $p$ can be determined for each province. This will be particularly advantageous if the ratios of the effective sample sizes of the panels are very different among the provinces; this is the case in SLID. In this connection, note that the effective sample sizes will be different between a province and the domain of movers from other provinces into that province. Given the relatively very small size of the domain of the interprovincial movers, a provincial weight adjustment factor (as in (6)) that would only account for the size of this domain but not for its different design effect would remain nearly optimal.

### 3.2 Non-identifiable temporal sample domains

It has been assumed thus far that the units of the non-overlap sample domain $s_a (\subset s_A)$ can be identified. However, the information needed to determine whether a unit in $s_A$ belongs to the frame domain $a$, of new entrants into the population after the start of the previous panel, may not be available for all units of $s_A$. In that situation the weighting process described above would combine the two samples $s_B$ and $s_A$ without distinguishing between the domains $s_{ab}$ and $s_a$ of $s_A$, so that the weights of units in $s_a$ would also be multiplied by $1-p$. The estimator $\hat{Y}_A^p$ in (3) would collapse then to

$$\hat{Y}_A^p = p\hat{Y}_{s_B} + (1-p)\hat{Y}_{s_A}. \tag{7}$$

The effect of this error is the underestimation of the total $Y_a$ for the population domain $a$ by the factor $p$. Part of the domain $a$, though, consists of newborns, which can be identified in $s_A$ with certainty. Their weights could very well be excepted from the adjustment by the factor $1-p$, but that would have no effect on cross-sectional estimation, unless newborns were part of the population of interest. Besides, adjusting the weights of newborns in $s_a$ by the factor $1-p$ has the desirable effect of producing a common household weight. A calibration of the weights of the combined sample to known population totals of the complete frame $A$ will lessen the under-representation of the rest of the domain $a$, which consists mainly of immigrants, but some bias may still result if the survey characteristics of the members of this part of the population are quite different from those of the members of the population domain $B$. Unless the time lag between the selection of the two panels is quite large, the size of this part of the population is very small, relative to the total population, and the potential bias effect on overall estimates of totals should be negligible.

The optimal (i.e., variance minimizing) value of $p$ in (7) is given now by

$$p'' = \frac{V(\hat{Y}_{s_A})}{V(\hat{Y}_{s_A}) + V(\hat{Y}_{s_B})} . \tag{8}$$

Disregarding finite population corrections it can be shown that (8) can be expressed as

$$\hat{p}_c'' = \frac{n_B d_A N_A^2 S_A^2}{n_B d_A N_A^2 S_A^2 + n_A d_B N_B^2 S_B^2} = \frac{n_B d_A}{n_B d_A + c n_A d_B} , \tag{9}$$

with $c = (N_B^2 S_B^2)(N_A^2 S_A^2)^{-1}$, and where $n_B$, $n_A$ are the sizes of the samples $s_B$ and $s_A$; $d_B$, $d_A$ are the design effects associated with $s_B$ and $s_A$ and the characteristic $y$; $N_A$, $N_B$ are the sizes of the frames $A$ and $B$; $S_A^2$, $S_B^2$ are the variances of the characteristic $y$ in $A$ and $B$. Noting that $N_B$ may be only a little smaller than $N_A$ (depending on the time lag between the two panels), and assuming that the unknown variances $S_A^2$ and $S_B^2$ are nearly equal, a good practical approximation of the optimal $p$ can be obtained by simply setting $c$ equal to one in (9). The assumption that the variances $S_A^2$ and $S_B^2$ are nearly equal is reasonable considering the magnitude of $N_B$ relative to that of $N_A$. Approximate values of $d_B$ and $d_A$ available from other surveys with the same designs as the two panels could be used, preferably for a characteristic such as the size of a large population domain. Now, if $\hat{Y}_c$ and $\hat{Y}_1$ denote the estimator $\hat{Y}_A^p$ in (7) when the weight adjustment $\hat{p}_c''$ in (9) is used with the true value of $c$ and the approximate value $c = 1$, respectively, then, ignoring finite population corrections, the loss of efficiency of $\hat{Y}_1$ relative to $\hat{Y}_c$ can be readily shown to be

$$\frac{Var(\hat{Y}_c) - Var(\hat{Y}_1)}{Var(\hat{Y}_c)} = -\frac{(c-1)^2}{c} p_1''(1 - p_1'') .$$

With a value of $c$ most likely in the neighbourhood of 1.0, the loss of efficiency will be negligible.

It is interesting to examine the efficiency of the estimator given by (7), with $p''$ as in (8), relative to the optimal estimator given by (3), with $p$ as in (4), used when the domain $s_a$ is identifiable. Let $\hat{Y}_A''$ and $\hat{Y}_A$ denote these estimators, respectively. Then, using the inequality $Cov^2(\hat{Y}_{s_A}, \hat{Y}_{s_{ab}}) \le Var(\hat{Y}_{s_A})Var(\hat{Y}_{s_{ab}})$ it can be shown that $Var(\hat{Y}_A) - Var(\hat{Y}_A'') \ge (p'' - p')Var(\hat{Y}_{s_A})$, where $p'$ is as in (5). If, as most likely, $Cov(\hat{Y}_{s_{ab}}, \hat{Y}_{s_a}) > 0$, then $p'' > p'$ and hence $Var(\hat{Y}_A) \ge Var(\hat{Y}_A'')$. Therefore, notwithstanding the use of the exact value of $p$ in the comparison, the approach taken in this subsection may actually result in reduction of the variance of derived estimators. A lower bound for the gain in efficiency would then be given by

$$\frac{Var(\hat{Y}_A) - Var(\hat{Y}_A'')}{Var(\hat{Y}_A'')} \ge \frac{(p'' - p')}{1 - p''} .$$

An extension of the weight adjustment procedure described above to surveys involving more than two panels with non-identifiable temporal sample domains is straightforward. There will be then as many weight adjustment factors, adding up to one, as there are panels. This very practical procedure will produce good cross-sectional estimates in multiple panel surveys in which the time lag between the selection of the panels is not large; otherwise the potential for bias due to the domain identification error may be of concern, mainly for estimates related to subpopulations composed in substantial proportion of new entrants.

-10-

## 4. THE WEIGHT SHARE METHOD FOR THE COMBINED PANELS

The combination of multiple panels for cross-sectional purposes involves an adjustment of the weights of the sampled units in the separate panels, as described above. Further weight adjustments are necessary because of the changes in the composition of panels after their selection. This section describes the application of a weight adjustment method, known as *the weight share method*, to the combined panel sample at any wave after the start of the most recent panel. Other weight adjustments are discussed in the next section.

The weight share method is a cross-sectional weighting procedure that assigns a basic weight to every individual in a panel at any wave after the first. In particular, the weight share method, as applied to a single panel, assigns a positive weight to non-selected individuals who join households containing at least one individual selected for the original sample. Following Lavallée (1995), in this paper such households are termed longitudinal households, while the non-selected individuals living in longitudinal households are termed cohabitants. The cohabitants are distinguished into originally present cohabitants if they belong to the original (sampled) population, and originally absent cohabitants if they are new entrants to the population. Other problematic situations that can be handled by the weight share method involve non-selected households formed after the first wave by members of different originally selected households, as well as originally selected individuals who have subsequently moved to other longitudinal households. For a detailed discussion of the weight share method, see Kalton and Brick (1995), and Lavallée (1995). Briefly outlined here, the weight share method for a single panel works as follows.

At a survey wave after the selection of the panel all individuals (including new members) in household $\mathcal{H}_i$, say, share a common weight defined as

$$w_i = \frac{1}{M_i} \sum_{k=1}^{M_i} w_{ik} \,,$$

where $M_i$ is the number of individuals in the household, at the time, that belong to the original (sampled) population, and $w_{ik}$ is the weight of the $k$-th individual in the household. For the individuals of the original population the weights are defined as random variables that take the value of the inverse of the inclusion probability if the individuals are in the original sample, and the value of zero otherwise, whereas for individuals not of the original population the weights are defined to be equal to zero. Then $E(w_{ik})=1$ for members of the original population, and hence $E(w_i)=1$ if $M_i \neq 0$, whereas $E(w_i)=0$ if $M_i=0$. Thus, every individual in household $\mathcal{H}_i$, for which $M_i \neq 0$, receives a weight that has expectation equal to one. Clearly then, the weight share method produces unbiased cross-sectional estimators of totals for the population of individuals living in households that contain at least one member of the original population, provided that the originally absent cohabitants (if any) can be identified for the correct specification of $M_i$. If the weights of the responding individuals at the time are adjusted for non-response, the relationship $E(w_i)=1$ may hold only approximately, and thus the resulting estimators may be only approximately unbiased.

For the purpose of applying the weight share method to a multiple panel survey the following concepts need to be considered. In multiple panel surveys, the original population for the combined panels is the union of the populations covered by the different panels at the time of their selection. Accordingly, the original sample consists of all selected units in the combined panel sample. Thus, originally present cohabitant is the individual that was eligible, but not selected, in any of the panels, or in a top-up sample. In this approach then, at any wave after the selection of the most recent panel a cohabitant is distinguished into originally present or originally absent with respect to the original combined panel sample, not with respect to each original panel. Notably, at the first wave of a new panel, or when a top-up sample is used, all cohabitants

-11-

are originally present. On the other hand, application of the weight share method separately to each panel (before combination) would require more precise information on the eligibility of the cohabitants for selection in each of the various panels, in order to distinguish into originally present and originally absent cohabitants and to identify the temporal domain that includes each of the cohabitants. Such information would most likely be unavailable. Moreover, combination of the panels after the weight share procedure would require a very complicated set of specifications to ensure that suitable weight adjustment factors are applied to each sampled unit. For instance, with the inclusion of the cohabitants into the panels through the weight share procedure, the frames of the panels will be different at each survey wave, thereby complicating the determination of the various temporal domains. It should also be pointed out that in multiple panel surveys sampled individuals may move from one panel to another panel between waves during the time period of panel overlap, and non-sampled households may be formed by members of originally selected households from different panels. Thus, the panels are truly distinct (and independent) only with respect to the time of their selection.

It follows from the forgoing considerations that the weight share method for multiple panels is to be applied to the combined panel sample, and not to each panel separately. Then, with the prescribed distinction of the two types of cohabitants, the case of the weight share method for a multiple panel survey reduces, essentially, to the case of a single panel survey. As a desirable consequence, the application of the weight share method to the combined sample yields always a common weight for all members of the same household. Following is an exemplification of the suggested weight share procedure for multiple panel surveys, involving the simple case of two panels.

Starting with a survey setting as depicted in Figure 1, with two overlapping panels at the time point of the start of the second panel, let there be $N$ individuals in the population at a later wave (time $t$), with $N_i$ individuals in household $\mathcal{H}_i$ ($i=1,...,H$) and $\sum N_i = N$. Let $M_i$ denote the number of individuals in household $\mathcal{H}_i$ that belong to the original population, with $M_{Bi}$ and $M_{ai}$ individuals from the domains $B$ and $a$, respectively, so that $M_i = M_{Bi} + M_{ai}$. Some, but not all, of the numbers $M_{Bi}$, $M_{ai}$ and $N_i - M_i$ may be zero for any particular household. Now, with the random weights of individuals in $B$ and $a$ as defined in Section 3.1, and with the weights of the $N_i - M_i$ originally absent cohabitants in $\mathcal{H}_i$ being equal to zero, the weight share method defines a common weight for any individual in $\mathcal{H}_i$ as

$$w_i = \frac{1}{M_i}\left[\sum_{k=1}^{M_{Bi}} w_{ik} + \sum_{k=1}^{M_{ai}} w_{ik}\right],$$ (10)

so that

$$E(w_t) = \frac{1}{M_t}\left[M_{Bt} + M_{ai}\right] = 1,$$

for each household for which $M_i \neq 0$. For the survey characteristic $y$, the total for the population of individuals at time $t$ can be expressed as $Y = \sum_{i=1}^{H} \sum_{k=1}^{N_i} y_{ik}$, where $y_{ik}$ is the value of $y$ for individual $k$ in household $\mathcal{H}_i$. Then, an estimator of $Y$ is given by

$$\hat{Y} = \sum_{i=1}^{H} w_i \sum_{k=1}^{N_i} y_{ik} = \sum_{i=1}^{H} w_i \left[\sum_{k=1}^{M_{Bi}} y_{ik} + \sum_{k=1}^{M_{ai}} y_{ik} + \sum_{k=1}^{N_i - M_i} y_{ik}\right] \doteq \hat{Y}_B + \hat{Y}_a + \hat{Y}_{A^c},$$ (11)

with $w_t$ as in (10), with $A^c$ denoting the set of individuals not in frame $A$, and with the obvious notation for the right hand side of (11). The estimator $\hat{Y}$ in (11) is given as the sum of three estimators, $\hat{Y}_B$, $\hat{Y}_a$ and $\hat{Y}_{A^c}$,

for the totals related to the population domains $B$, $a$ and $A^c$, respectively. The estimators $\hat{Y}_B$ and $\hat{Y}_a$ are unbiased, even though they are based on sets of units that may not be identical to the original sample domains $s_B \cup s_{ab}$ and $s_a$, respectively. For example, the estimator $\hat{Y}_B$ is based on a set of units consisting of the remaining units of the original sample domain $s_B \cup s_{ab}$ from frame $B$, and possibly of cohabitants originally present in $B$. The estimator $\hat{Y}_{A^c}$ is not unbiased for $Y_{A^c}$, because individuals in $A^c$ who leave in households that contain no members of the original population are not represented in the panel survey. Nevertheless, the estimator $\hat{Y}_{A^c}$ is unbiased for the total corresponding to the rest of $A^c$, which is represented in the combined panels by the originally absent cohabitants.

In the special case when time $t$ coincides with the start of the second panel (or with the time of selection of a top-up sample), $A^c = \emptyset$, $N_t = M_t$, and the estimator $\hat{Y} = \hat{Y}_B + \hat{Y}_a$ is unbiased for $Y$. This strongly argues in favour of using a top-up sample (or a supplementary sample from $A^c$) at each wave of a panel survey for cross-sectional purposes.

The situation in which the members of the population domain $a$ cannot be identified in the sample deserves special attention. Then, following the discussion in Section 3.2, for a household $\mathcal{H}_i$ containing members from the domain $a$ the weight share method produces a household weight $w_i$ with

$$E(w_i) = \frac{1}{M_i}\left[M_{Bi} + (1-p)M_{ai}\right] = 1 - p\frac{M_{ai}}{M_i} \ . \tag{12}$$

The expected value of the estimator $\hat{Y}$ is then given by

$$E(\hat{Y}) = \sum_{i=1}^{H}\sum_{k=1}^{N_i} y_{ik} - p\sum_{i=1}^{H}\frac{M_{ai}}{M_i}\left[\sum_{k=1}^{M_{Bi}} y_{ik} + \sum_{k=1}^{M_{ai}} y_{ik} + \sum_{k=1}^{N_i - M_i} y_{ik}\right]. \tag{13}$$

For households with no members from $a$ (i.e., with $M_{ai} = 0$) the bias term in the right-hand side of equation (13) is zero. For households composed solely of individuals from $a$ (i.e., with $N_i = M_{ai}$) it follows that $E(w_i) = 1 - p$, for all household members. In this case the effect of wrongly adjusting the initial weights of members of the domain $a$ on estimation is as before the application of the weight share method; see discussion in Section 3.2. It also follows from (13) that when $0 < M_{ai} < M_i$ all $N_i$ household members will be affected by the weight adjustment error if the household contains originally absent cohabitants, or if some household members from one or both of the domains $B$ and $a$ are (selected) movers into the household, unless all $M_{ai}$ household members are originally present cohabitants. The number of such households in the combined sample is likely to be very small, and hence the spread of the identification error beyond the domain $a$ will not be of serious consequence.

As with the weight adjustment involved in the combination of panels, the weight share adjustment may also be carried out at a superstratum level, say province, for the combined sample of each province. In this approach, those individuals who at time $t$ reside in a province other than the one in which they resided at the time of selection of any of the panels are treated as originally absent, since they were not members of the original population of their new province. In particular, interprovincial movers (selected or non-selected in their original province) who are found in longitudinal households in their new province at time $t$ are treated as originally absent cohabitants. Households made up solely of selected movers from another province are then discarded. The part of the population represented by the discarded movers is very small, and any implied bias effect will be lessened by a calibration of the sample weights to known current population totals. This part of the population is covered when a new panel (or a top-up sample) is selected, and no bias is then incurred, provided that its members can be identified in the new panel or in the top-up sample as new entrants

into their current province, so that they can be excluded from the weight adjustment that combines the panels. It is to be noted that when a top-up sample is used, interprovincial movers (selected or non-selected) who are found in longitudinal households in their new province at time $t$ are treated in this approach as originally present cohabitants. Obviously, similar would be the treatment of individuals in a panel who have moved to another province before the selection of the next panel. In an empirical study based on data from the third wave of the first panel of SLID, it was found that there is negligible loss of efficiency due to discarding of interprovincial movers of the aforementioned type for most of the studied survey characteristics. For some characteristics there is in fact gain in efficiency, as measured by difference in CV's, primarily because this alternative weight share procedure avoids the inflationary effect on variances that is associated with large differences in magnitude between the weights of the interprovincial movers and the weights of original members of the movers' new province. The application of the weight share procedure separately for each province enjoys certain operational advantages over the standard weight share procedure. An account of the comparative merits of the two approaches is given in Merkouris (1999).

Finally, it is important to note that the estimator $\hat{Y}$ in (11) can be expressed as

$$\hat{Y} = \sum_{t=1}^{H} w_t Y_t \, ,$$

where $Y_t = \sum_{k=1}^{N_t} y_{tk}$ is the total for household $\mathcal{H}_t$. Thus, $\hat{Y}$ is also an estimator of the household level total at time $t$. When time $t$ coincides with the start of the second panel, or when the survey employs a top-up sample, the estimator $\hat{Y}$ is unbiased.


## 5. INTEGRATION OF WEIGHT ADJUSTMENTS

In addition to the weight adjustments described so far, other adjustments to the weights of a panel household survey may also be required. The integration of the various weight adjustments is briefly outlined below.

The first adjustment, applied in relation to the original sample units, is for wave non-response, which arises when a sampled unit responds for some but not all of the waves for which it was eligible. For a discussion on weight adjustment for wave non-response, see Kalton and Brick (1995). The adjustment is made separately to the different panels at each wave.

The second adjustment is for the combination of the samples of the various panels into one sample for cross-sectional estimation. It applies to the weights of the sampled units of the panels, adjusted for wave non-response, and employs the method described in Section 3.

The third adjustment involves the application of the weight share procedure to the combined panel sample at any wave after the start of the most recent panel, as described in Section 4.

Finally, in the weight calibration adjustment the weights of the combined panel units are adjusted so as to make the estimated totals of certain auxiliary characteristics equal to known population totals for these characteristics at the current wave, which in the simple case as in Figure 1 correspond to totals of the complete frame $A$. In more general situations, after the selection of the most recent panel the calibration totals will include the new entrants into the population. Note that in the absence of a top-up sample the new entrants will be represented in the panels only by the originally absent cohabitants. Separate calibrations of the weights of the combined sample for each of the different temporal domains (when the panel units from these domains can be identified) to corresponding population totals may not be feasible or sensible for reasons already noted in Section 3.

## 6. CONCLUDING REMARKS

The weighting procedures described in this paper can be used to combine information from multiple panels of a repeated household survey for cross-sectional estimation in a fairly general setting involving panels with given designs. Design issues regarding determination of optimal sampling fractions for the panels, in conjunction with efficient combination of the panel data, are beyond the scope of this paper. The proposed estimation procedures are operationally convenient for any number of overlapping panels, and for different situations regarding the identifiability of various temporal panel domains. It has been shown that although a multiple panel survey can be viewed as a special type of multiple frame survey, its distinctive dynamic character renders conventional multiple frame estimation procedures problematic or even not applicable. Theoretical and practical issues related to the application of a weight share adjustment, to the calibration weight adjustment and to the integration of the various weighting procedures involved in a multiple panel survey have also been addressed. A detailed empirical study of issues pertaining to the determination of weight adjustment factors for combining two panels of SLID, based on the methodology of this paper, is described in Latouche et al. (1999). The variance of cross-sectional estimators has been discussed in this paper only in the context of efficient combination of panels. Variance estimation issues related to changes in the sample over time, particularly to moves from stratum to stratum, are discussed in Merkouris (1999). A jackknife variance estimation procedure for surveys with multiple frames (possibly sampled with different designs) discussed in Lohr and Rao (1997) is applicable in the context of multiple panel surveys. It is to be remarked, in conclusion, that the quality of a cross-sectional estimation procedure depends on the identifiability of various overlap temporal sample domains; on design features of the survey, such as the duration of (and the lag between) the panels, or the use of a supplementary sample at any survey wave; on the adequacy of the information on cohabitants required for the application of the weight share method.

## ACKNOWLEDGEMENTS

## REFERENCES

DEVILLE, J. C. (1998). Les enquetes par panel : En quoi different-elles des autres enquetes? Suivi de comment attraper une population en se servant d'une autre. *Proccedings of the Journées Francofones sur les sondages, Rennes, France* (to appear).

KALTON, G., and ANDERSON, D.W. (1986). Sampling rare populations. *Journal of the Royal Statistical Society* Ser.A, 149, 65-82.

KALTON, G., and CITRO, C.F. (1993). Panel Surveys: Adding the fourth dimension. *Survey Methodology*, 19, 205-215.

KALTON, G.,and BRICK, J.M.(1995). Weighting schemes for household panel surveys. *Survey Methodology*, 21, 33-44.

LATOUCHE, M., DUFOUR, J., and MERKOURIS, T. (1999). SLID cross-sectional weighting: Combining two or more panels. Internal document, Statistics Canada.

LAVALLÉE, P.(1994). Ajout du second panel à l'EDTR: sélection et pondération. Internal document. Statistics Canada.

LAVALLÉE, P. (1995). Cross-sectional weighting of longitudinal surveys of individuals and households using the weight share method. *Survey Methodology*, 21, 25-32.

LAVIGNE, M., and MICHAUD, S. (1998). General aspects of the Survey of Labour and Income Dynamics. Working Paper SLID 98-05 E, Statistics Canada.

LOHR, S., and RAO, J.N.K. (1997). Jackknife variance estimation in multiple frame surveys. *Proceedings of the Section on Survey Research Methods, Americal Statistical Association*, 552-557.

MERKOURIS, T. (1999). The weight share method for panel household surveys: Issues related to moves between strata. Methodology Branch Working Paper HSMD 99-003E, Statistics Canada.

SINGH, A.C., and WU, S. (1996). Estimation for multiframe complex surveys by modified regression. *Proceedings of the Survey Methods Section, Statistical Society of Canada*, 69-77.

SKINNER, C.J., and RAO, J.N.K. (1996). Estimation in dual frame surveys with complex designs. *Journal of the American Statistical Association*, 91, 349-356.