

Register-based sampling for household panels

Citation for published version (APA):

van den Brakel, J. (2016). Register-based sampling for household panels. *Survey Methodology*, 42(1), 137-159.

Document status and date:

Published: 01/01/2016

Document Version:

Publisher's PDF, also known as Version of record

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.umlib.nl/taverne-license

Take down policy

If you believe that this document breaches copyright please contact us at:

repository@maastrichtuniversity.nl

providing details and we will investigate your claim.

Catalogue no. 12-001-X
ISSN 1492-0921

Survey Methodology

Register-based sampling for household panels

by Jan A. van den Brakel

Release date: June 22, 2016



Statistics
Canada

Statistique
Canada

Canada

How to obtain more information

For information about this product or the wide range of services and data available from Statistics Canada, visit our website, www.statcan.gc.ca.

You can also contact us by

email at STATCAN.infostats-infostats.STATCAN@canada.ca

telephone, from Monday to Friday, 8:30 a.m. to 4:30 p.m., at the following toll-free numbers:

- Statistical Information Service 1-800-263-1136
- National telecommunications device for the hearing impaired 1-800-363-7629
- Fax line 1-877-287-4369

Depository Services Program

- Inquiries line 1-800-635-7943
- Fax line 1-800-565-7757

Standards of service to the public

Statistics Canada is committed to serving its clients in a prompt, reliable and courteous manner. To this end, Statistics Canada has developed standards of service that its employees observe. To obtain a copy of these service standards, please contact Statistics Canada toll-free at 1-800-263-1136. The service standards are also published on www.statcan.gc.ca under “Contact us” > “Standards of service to the public.”

Note of appreciation

Canada owes the success of its statistical system to a long-standing partnership between Statistics Canada, the citizens of Canada, its businesses, governments and other institutions. Accurate and timely statistical information could not be produced without their continued co-operation and goodwill.

Standard table symbols

The following symbols are used in Statistics Canada publications:

- . not available for any reference period
- .. not available for a specific reference period
- ... not applicable
- 0 true zero or a value rounded to zero
- 0^s value rounded to 0 (zero) where there is a meaningful distinction between true zero and the value that was rounded
- ^P preliminary
- ^r revised
- X suppressed to meet the confidentiality requirements of the *Statistics Act*
- ^E use with caution
- F too unreliable to be published
- * significantly different from reference category ($p < 0.05$)

Published by authority of the Minister responsible for Statistics Canada

© Minister of Industry, 2016

All rights reserved. Use of this publication is governed by the Statistics Canada [Open Licence Agreement](#).

An HTML version is also available.

Cette publication est aussi disponible en français.

Register-based sampling for household panels

Jan A. van den Brakel¹

Abstract

In the Netherlands, statistical information about income and wealth is based on two large scale household panels that are completely derived from administrative data. A problem with using households as sampling units in the sample design of panels is the instability of these units over time. Changes in the household composition affect the inclusion probabilities required for design-based and model-assisted inference procedures. Such problems are circumvented in the two aforementioned household panels by sampling persons, who are followed over time. At each period the household members of these sampled persons are included in the sample. This is equivalent to sampling with probabilities proportional to household size where households can be selected more than once but with a maximum equal to the number of household members. In this paper properties of this sample design are described and contrasted with the Generalized Weight Share method for indirect sampling (Lavallée 1995, 2007). Methods are illustrated with an application to the Dutch Regional Income Survey.

Key Words: Probabilities proportional to size; Indirect sampling; Consistent weighting of persons and households; Regional Income Survey; Generalized Weight Share method.

1 Introduction

Statistics Netherlands conducts two important sample surveys to describe the income and wealth situation of the Dutch population. First, the Dutch Regional Income Survey (RIS) provides a description of the income and wealth situation, being accurate at a very detailed regional level. This is accomplished by publishing accurate income distributions for persons and households at a level of neighbourhoods on a yearly basis, using a large sample based on a small set of the main income components derived in a relatively straightforward manner from tax administration. Second, the Income Panel Survey (IPS) publishes yearly income and wealth characteristics of the Dutch population at a more aggregated regional level. This survey is based on a large set of variables using all possible income components of households that can be derived from the available administrative data in the Netherlands. The derivation of the variables for this survey is more time consuming. Therefore the sample size of this survey is considerably smaller than the RIS. Both surveys are designed as a household panel where both person and household based variables about income and wealth are observed.

Households are often considered as the sampling units in panels conducted to collect information at the level of households and persons (Lynn 2009; Smith, Lynn and Elliot 2009). Such panels are used for longitudinal analysis as well as the production of cross-sectional estimates. Using households as the sampling units in a panel design has, however, some major disadvantages due to their instability over time. As time proceeds, households might disintegrate, join or split, new members might enter the households and other members might leave the households for different reasons. Kalton and Brick (1995) explain that these changes can affect the selection probabilities of the households in the sample. Reconstruction of the correct inclusion probabilities of the sampling units is essential to derive correct weights for analysis purposes, in particular if the panel is used for producing cross-sectional estimates.

1. Jan A. van den Brakel, Department of Statistical Methods, Statistics Netherlands, P.O. Box 4481, 6401 CZ Heerlen, The Netherlands and Department of Quantitative Economics, Maastricht University School of Business and Economics, P.O. Box 616, 6200 MD, Maastricht, The Netherlands. E-mail: ja.vandenbrakel@cbs.nl.

Consider a panel where households are selected by means of simple random sampling, say at time $t = 0$. In many panels, people that enter a sampled household at a later stage are also included in the panel. These individuals are called cohabitants by Lavallée (1995). As time proceeds, more and more cohabitants are included in the sample and disturb the equal probability design that is used to select the initial sample (Kalton and Brick 1995). Consider for example household A, which is selected in the sample when the panel started at $t = 0$. If after some period of time this household merges with another household B, which was initially not selected for the panel at time $t = 0$, then the selection probability of this new household is the sum of the selection probabilities of households A and B at time $t = 0$. Not correcting for differences in selection probabilities due to the gradual increasing share of cohabitants in the sample leads to biased inference. Ernst (1989) proposes the Weight Share method to overcome these problems. Lavallée (1995) extends this method to the Generalized Weight Share method as a solution for drawing inference about target populations that are sampled through the use of a frame that refers to a different population.

The RIS and the IPS are both based on a panel and are conducted to collect information about households and persons. To avoid the problems with panels using households as sampling units, an alternative design is applied. Instead of households, so-called core persons are drawn with an equal probability design, who are followed over time. All household members belonging to the household of a core person at each particular period are included in the sample. This results in a sample design where households are drawn proportionally to the household size and households can be selected more than once, but with a maximum that is equal to the household size. This design is an application of indirect sampling (Lavallée 1995, 2007; Deville and Lavallée 2006).

The purpose of this paper is to describe a sample design with an estimation technique that is useful for panels that collect information at person and household level. The methodology employed in this paper is particularly useful for register based sampling, since the core persons are included in the sample indefinitely. The sample design is also useful for Web panels, but might require some form of rotating design to avoid problems with panel attrition. This means that sampling units enter the panel, are observed multiple times and leave the panel according to a pre-specified pattern (Smith et al. 2009). The main contribution of this paper to the existing literature is that explicit expressions for the variance of the target parameters are derived using inclusion expectations instead of inclusion probabilities under the aforementioned sample design. A measure of the minimum accuracy for an estimated income distribution is proposed and explicit expressions for the minimum sample size are derived. The RIS is used throughout the paper to illustrate the described sampling techniques.

The paper is organized as follows. A description of the sample design of the RIS is given in Section 2. In Section 3 the concept of inclusion expectations is introduced as a convenient practical alternative for inclusion probabilities. Subsequently, first and second order inclusion expectations are derived for the proposed sampling design. These inclusion expectations are required to construct the π -estimator or Horvitz-Thompson (HT) estimator (Narain 1951; Horvitz and Thompson 1952). It is also shown that the same weights can be derived as a special case of the Generalized Weight Share method for indirect sampling (Lavallée 1995, 2007). The key target variables for the RIS are estimated income distributions. In Section 4 formulas for the minimum required sample size are derived based on a precision measure for estimated income distributions. Since households can be selected more than once, an expression for the expected number of unique households is derived in Section 4. The estimation procedure of the RIS is based on linear weighting using the general regression (GREG) estimator (Särndal, Swensson and Wretman 1992) and is

described in Section 5. The integrated weighting method of Lemaître and Dufour (1987), Nieuwenbroek (1993) and Steel and Clark (2007) is applied to obtain equal weights for persons belonging to the same household. In Section 6 variance approximations for the GREG estimator under the proposed sample design are derived. An application to the RIS is provided in Section 7. The paper concludes with a discussion in Section 8.

2 Sampling design

The target population of the RIS is all natural persons residing in the Netherlands. The sample frame is a register containing all natural persons aged 15 years and over residing in the Netherlands as far as they are known to the Tax Office. From this register a stratified simple random sample of so-called core persons is drawn with a sample fraction of 0.16. Neighbourhoods are used as the stratification variable. Although an equal probability design is used, stratified sampling is useful to eliminate the variation between strata and to meet minimum precision requirements for the individual strata. The Netherlands is divided in about 2,830 neighbourhoods with an average size of 5,000 persons aged 15 years and over.

The RIS has been conducted as a panel since 1994. A first requirement for correct cross-sectional inference with this panel is to have correct first and second order inclusion expectations for the sampling units, which are derived in Section 3. A second requirement for correct cross-sectional inference is to keep the panel representative of the target population. To this end, it is determined on a yearly basis which part of the population has entered the target population of the RIS through birth and immigration. From this subpopulation, a stratified simple random sample of core persons with a sample fraction of 0.16 is selected. These core persons are added to the panel of the RIS, with the purpose to maintain a representative sample.

Neighbourhoods are the most detailed level of publication for the RIS and are therefore used as strata. In Section 4 expressions for minimum sample sizes based on precision requirements are derived. Core persons remain in the panel indefinitely. On each survey occasion, all members of the core person's household are also included in the sample. Persons that leave the household of a core person also leave the panel. New persons entering the household of the core person are followed in the panel as long as this person stays in the household of a core person. Information about the household composition of the core persons are obtained from the Municipal Basis Administration (MBA), which is the Dutch government's registry of all residents in the country. Dutch citizens are required by law to report changes in their demographics to their municipalities. The MBA is used in combination with the information from tax administrations to identify household members of the core persons in the sample.

The sample design results in a sample of households where the households are selected with probabilities proportional to the number of persons aged 15 years or older belonging to a household at the current period. Households can be selected more than once, but with a maximum that equals the number of household members aged 15 year or older. In this paper the term core persons is used to refer to the persons that are initially included in the sample and are followed over time in the panel. The term persons is used to refer to the sample obtained if all the household members at a particular period are included in the sample.

The IPS applies a similar sample design with a substantially smaller sampling fraction. The RIS, like the IPS, are register based samples which implies that for each person that is included in the sample, the necessary information for the RIS variables is obtained from the registers of the Tax Office. Core persons

and their household members are therefore not aware that they are included in these samples. This has the advantage that there are no problems with selective non-response and panel attrition. This also makes it possible to include the core persons indefinitely. In the case of a panel where sampling units must complete a questionnaire, some kind of rotating design would be required in order to avoid selectivity bias due to panel attrition. Also, problems with measurement bias associated with data collection where sampling units are asked to complete a questionnaire do not occur. Of course other types of measurement errors are encountered with a survey that is based on registrations (Wallgren and Wallgren 2007). It is assumed that all the required information about income to estimate the target parameters of the RIS and the IPS are available in these registers. Since all the required information is available in a register, a complete enumeration of the population is possible. In the past, however, the IT infrastructure was insufficient to produce timely regional income statistics based on a complete enumeration of the Dutch population. Therefore the RIS was traditionally based on a large sample with a fraction of 0.16 core persons. For the same reason the IPS is traditionally based on a sample of about 80,000 core persons. With the current computational capacity a complete enumeration would still be very demanding but not impossible. The main rationale for conducting this survey as a sample is to maintain the panel for longitudinal analysis that cover time periods from the past where a census was impossible.

3 Inclusion weights

3.1 Weighting with inclusion expectations

For design-based inference, first and second order inclusion probabilities for households and persons are required. Let M denote the number of households in the population, N the number of persons in the population aged 15 years or over and g_k the number of persons aged 15 years or over that belong to the k^{th} household. With the sample design described in Section 2, households k can be included more than once but a maximum of g_k times. This complicates the derivation of inclusion probabilities since the probability of selecting household k is equal to the selection probability of the union of its household members (k, j) aged 15 years and over. This probability is defined as:

$$\begin{aligned}
 P(k \in s) &= P\left(\bigcup_{j=1}^{g_k} [(k, j) \in s]\right) = \sum_{j=1}^{g_k} P((k, j) \in s) \\
 &\quad - \sum_{j=1}^{g_k} \sum_{j'=j+1}^{g_k} P([(k, j) \cap (k, j')] \in s) \\
 &\quad + \sum_{j=1}^{g_k} \sum_{j'=j+1}^{g_k} \sum_{j''=j'+1}^{g_k} P([(k, j) \cap (k, j') \cap (k, j'')] \in s) - \dots
 \end{aligned}$$

This kind of computation can be avoided by using the concept of inclusion expectations instead of inclusion probabilities. Bethlehem (2009), Chapter 2, generalizes the HT estimator to the concept of inclusion expectation for sampling with replacement. Let a_k denote the number of times that household k is selected in the sample. In the proposed sample design $a_k \in [0, 1, \dots, g_k]$. Let $E(\cdot)$ denote the expectation with

respect to the sample design. Now $\pi_k = E(a_k)$ denotes the inclusion expectation of sampling unit k . Since a_k can be larger than one, π_k can also take values larger than one and can therefore no longer be interpreted as an inclusion probability. It can, however, be interpreted as an expectation.

The parameter of interest is the population total, which is defined as

$$t_y = \sum_{k=1}^M \sum_{j=1}^{N_k} y_{kj} \equiv \sum_{k=1}^M y_k. \tag{3.1}$$

The HT estimator for the population total in (3.1) can be defined as

$$\hat{t}_y = \sum_{k=1}^M \frac{a_k y_k}{\pi_k}. \tag{3.2}$$

Since $E(a_k) = \pi_k$, it follows that this HT estimator is design unbiased. Let $\pi_{kk'}$ denote the inclusion expectation of units k and k' , i.e., $\pi_{kk'} = E(a_k a_{k'})$. The variance of the HT estimator is by definition equal to

$$\begin{aligned} V(\hat{t}_y) &= \sum_{k=1}^M \sum_{k'=1}^M \text{Cov}(a_k a_{k'}) \frac{y_k}{\pi_k} \frac{y_{k'}}{\pi_{k'}} \\ &= \sum_{k=1}^M \sum_{k'=1}^M [E(a_k a_{k'}) - E(a_k) E(a_{k'})] \frac{y_k}{\pi_k} \frac{y_{k'}}{\pi_{k'}} \\ &= \sum_{k=1}^M \sum_{k'=1}^M (\pi_{kk'} - \pi_k \pi_{k'}) \frac{y_k}{\pi_k} \frac{y_{k'}}{\pi_{k'}}. \end{aligned}$$

Note that in the case of sampling without replacement a_k is a dummy taking values zero or one indicating whether unit k is selected in the sample. In this case π_k and $\pi_{kk'}$ are the usual first and second order inclusion probabilities. This illustrates that the standard HT estimator, based on inclusion probabilities, can be extended easily to inclusion expectations. In the case of sample designs where units can be selected more than once, it is more convenient to work with inclusion expectations, since they are derived relatively easily. In the remainder of this subsection, first and second order inclusion expectations for the sample design described in Section 2 are derived.

Core persons are drawn by means of stratified simple random sampling. Since stratification is based on geographical regions, all members of a household k belong to the same stratum h at the moment of drawing core persons. Let N_h denote the number of persons in the population of stratum h aged 15 years or over, n_h the number of core persons selected in the sample from stratum h and g_k the number of persons aged 15 years or over, belonging to household k . Finally, a_{jk} denotes an indicator that is equal to one if person j from household k is selected in the sample and zero otherwise. The first order inclusion expectation of the k^{th} household equals

$$\pi_{kh} = E(a_k) = E\left(\sum_{j=1}^{g_k} a_{jk}\right) = \sum_{j=1}^{g_k} E(a_{jk}) = g_k \frac{n_h}{N_h}. \tag{3.3}$$

Second order inclusion expectations for households k and k' for $k \neq k'$ belonging to the same stratum h , equal

$$\pi_{kk'} = E(a_k a_{k'}) = E\left(\sum_{j=1}^{g_k} a_{jk} \sum_{j'=1}^{g_{k'}} a_{j'k'}\right) = \sum_{j=1}^{g_k} \sum_{j'=1}^{g_{k'}} E(a_{jk} a_{j'k'}) = g_k g_{k'} \frac{n_h (n_h - 1)}{N_h (N_h - 1)}. \quad (3.4)$$

The second order inclusion expectation for household $k = k'$ from the same stratum h , is given by

$$\begin{aligned} \pi_{kk} &= E(a_k a_k) = E\left(\sum_{j=1}^{g_k} a_{jk} \sum_{j'=1}^{g_k} a_{j'k}\right) = E\left(\sum_{j=1}^{g_k} a_{jk} + \sum_{j=1}^{g_k} \sum_{j' \neq j=1}^{g_k} a_{jk} a_{j'k}\right) \\ &= \sum_{j=1}^{g_k} E(a_{jk}) + \sum_{j=1}^{g_k} \sum_{j' \neq j=1}^{g_k} E(a_{jk} a_{j'k}) = g_k \frac{n_h}{N_h} + g_k (g_k - 1) \frac{n_h (n_h - 1)}{N_h (N_h - 1)}. \end{aligned} \quad (3.5)$$

Second order inclusion expectations for households k and k' for $k \neq k'$ belonging to two different strata h and h' equal

$$\pi_{kk'} = E(a_k a_{k'}) = E\left(\sum_{j=1}^{g_k} a_{jk} \sum_{j'=1}^{g_{k'}} a_{j'k'}\right) = \sum_{j=1}^{g_k} \sum_{j'=1}^{g_{k'}} E(a_{jk} a_{j'k'}) = g_{kh} g_{k'h'} \frac{n_h n_{h'}}{N_h N_{h'}}. \quad (3.6)$$

An alternative proof based on the definition of an expected value, which does not use the rule that the expected value of a sum of mutual dependent variables is equal to the sum over the expected values of these variables is given by van den Brakel (2013).

As time proceeds the household composition of the core persons changes, which affects the inclusion expectations of the households in the sample. If sampling fractions differ between strata, the inclusion expectations (3.3) through (3.6) become more complicated and require information of stratum membership for all persons belonging to the household of the core persons. This complication is avoided by choosing a self-weighted sampling design. In this case each household member of a core persons has the same inclusion probability and the only household specific information required to derive household inclusion expectations is the number of persons aged 15 years and over in the household of the core person.

Since all members of a selected household are included in the sample, it follows that the first order inclusion expectations for persons belonging to household k are equal to the first order inclusion expectation of household k defined in (3.3). The second order inclusion expectations for persons from two different households k and k' , are equal to (3.4) for two households from the same stratum or (3.6) for two households from two different strata. The second order inclusion expectations for persons from the same household are defined by (3.5).

During the review the question was raised whether the inclusion expectations themselves have a variance that should be taken into account in the variance of HT or GREG estimators when they are based on inclusion expectations instead inclusion probabilities. In the finite population each person and each household has a pre-specified inclusion expectation. For the households observed in the sample these expectations can be calculated exactly without uncertainty since all information required to evaluate the true value of these expectations is available. Substituting inclusion probabilities for expectations, therefore does not result in an additional variance component.

3.2 Generalized Weight Share method

The sample design described in Section 2 can be considered as a special case of indirect sampling (Lavallée 2007). Indirect sampling refers to the situation where the population of interest is sampled through the use of a frame that refers to a different population. Lavallée (1995) develops the Generalized Weight Share method to construct weights for these situations and can be used to derive design weights for households and persons in the sample design described in Section 2.

Following the notation of Lavallée (1995) for the case of indirect sampling, there is a population U^A of size N^A from which a sample s^A of size n is drawn with selection probabilities π_i^A . In addition, there is the target population U^B of size N^B . This population can be divided in M^B clusters. Each cluster k contains N_k^B units, such that $N^B = \sum_{k=1}^{M^B} N_k^B$. The situation for the sample design described in Section 2 is depicted in Figure 3.1. The clusters are households, U^A is the population of persons aged 15 years and over, and U^B is the population of all persons residing in the Netherlands. Persons in U^A and U^B are depicted as circles, households in U^B are depicted as shaded squares, and the circles within a shaded square visualise persons belonging to the same household. Figure 3.1 shows respectively, a single person household, a two person household containing for example a divorced parent with a child younger than 15, a two person household containing two adults without children, and a four person household containing two parents with two children and one of the children is younger than 15 while the other is 15 years or older. The arrows depict the links between the units of U^A and U^B . In the sample design considered in Section 2, each unit in U^A has exactly one unique link with a unit in U^B . Clusters in U^B have at least one link with units in U^A . Links are identified with an indicator variable

$$l_{ij} = \begin{cases} 1 & \text{if there is a link between } i \in U^A \text{ and } j \in U^B \\ 0 & \text{if there is no link between } i \in U^A \text{ and } j \in U^B. \end{cases}$$

If a unit i in U^A is selected in the sample, the entire cluster k to which this unit belongs, is included in the sample. The parameter of interest is the population total in U^B and is similar to (3.1) defined as $t_y = \sum_{k=1}^{M^B} \sum_{j=1}^{N_k^B} y_{kj}$. An estimator for t_y is defined as

$$\hat{t}_y = \sum_{k=1}^m \sum_{j=1}^{N_k^B} w_{kj} y_{kj}, \quad (3.7)$$

with m the number of unique clusters (households) included in the sample and w_{kj} the weight attached to each unit j of cluster k . Generally the inverse of the selection probabilities of units (k, j) observed in the sample are used as weights in the HT estimator. In this situation not all units in the sample have a known inclusion probability. Firstly not all units in U^B have a link to U^A . Secondly, as time proceeds household compositions change due to marriages, divorces, departures of children and cohabitation. As a result, as time proceeds, units with a link to U^A enter the clusters in the sample although they are not initially included in the sample drawn from U^A . For these units inclusion probabilities are not necessarily known. They affect, however, the inclusion expectations of the clusters included in the sample. Reconstruction of the inclusion probabilities requires information of selection probabilities of all units in the population at the moment that the sample is drawn. In many practical situations this information is not available.

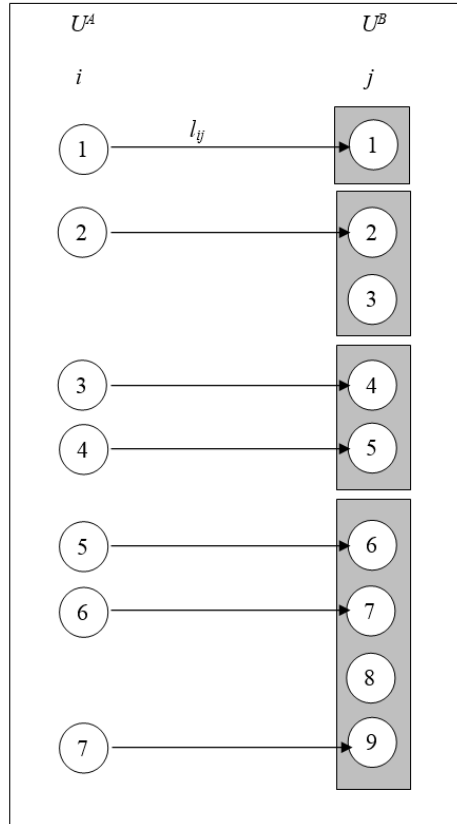


Figure 3.1 Links between units from the sample frame and units from the target population.

The Generalized Weight Share method can be used to derive non-zero weights for all units in the sample. This method starts by deriving initial weights, which are defined as

$$w_{kj}^* = \begin{cases} \frac{\delta_i^A}{\pi_i^A} & \text{if } (k, j) \text{ has a link with } i \in U^A \\ 0 & \text{otherwise} \end{cases},$$

with δ_i^A an indicator variable that is equal to one if i is included in the sample s^A and zero otherwise. This expression follows directly from Lavallée (1995), equation (2) in combination with the fact that in this application each unit in U^A has exactly one unique link with a unit in U^B , see Figure 3.1. In a second step a so-called basic weight for each cluster k is derived as the mean of all initial weights within each cluster

$$w_k = \frac{\sum_{j=1}^{N_k^B} w_{kj}^*}{\sum_{j=1}^{N_k^B} l_{kj}},$$

which follows from Lavallée (1995), equation (7). Finally all persons j that belong to the same household k receive the same weight assigned to their household, i.e., $w_{kj} = w_k$ for all $j \in k$. A proof that the use of the basic weights in (3.7) is an unbiased estimator for the population total is also given by Lavallée (1995).

Let $\sum_{j=1}^{N_k^B} I_{kj} = g_k$ denote the number of persons in household k aged 15 years and older and a_k the number of core persons in household k , i.e., the number of persons in household k that are included in sample s^A . Since s^A is drawn by means of stratified simple random sampling, it follows that $\pi_i^A = n_h^A / N_h^A$ with N_h^A the number of persons aged 15 years and older in the population of stratum h , and n_h^A the number of core persons selected in the sample from stratum h . Then it follows that

$$w_k = \frac{a_k}{g_k} \frac{N_h^A}{n_h^A}. \tag{3.8}$$

Inserting the first order inclusion expectation (3.3) into (3.2) gives the same HT estimator as derived with the Generalized Weight Share method, i.e., inserting (3.8) into (3.7).

The derivation of the inclusion expectations in Subsection 3.1 applies to stratified sampling of households with inclusion expectations proportional to household size and is a special case of the Generalized Weight Share method. An argument to apply a design as outlined in Section 2 is that sampling households proportional to household size is efficient for target variables that are positively correlated with household size.

Lavallée (1995) also provides variance expressions for (3.7) based on the Generalized Weight Share method. This expression is based on the first and second order inclusion probabilities of the sample units drawn from U^A and a transformation of the target variable. As a result the property that clusters are drawn proportional to their size is not made explicit, nor that the fact they are drawn partially with replacement. In Section 6 it is pointed out that the variance expressions in Lavallée (1995) for this application are equal to the variance expressions based on the inclusion expectations derived in (3.3) through (3.6).

4 Sample size determination

The purpose of the RIS is to publish income distributions for households and persons at different geographical levels. Income distributions for households for region or area r are defined as

$$P_{lr} = \frac{M_{lr}}{M_{+r}}, \quad l = 1, \dots, L, \tag{4.1}$$

where M_{lr} denotes the number of households from region r , belonging to the l^{th} income category, and $M_{+r} = \sum_l M_{lr}$, the total number of households in area r . This income distribution is estimated as

$$\hat{P}_{lr} = \frac{\hat{M}_{lr}}{M_{+r}}, \quad l = 1, \dots, L, \tag{4.2}$$

where \hat{M}_{lr} denotes an appropriate direct estimator for the total number of households from area r , classified to the l^{th} income category. For the moment the HT estimator is assumed as an appropriate estimator for M_{lr} , i.e.,

$$\hat{M}_{lr} = \sum_{h \in r} \sum_{k=1}^{m_h} \frac{y_{khl}}{\pi_k},$$

where $y_{khl} = 1$ if household k from stratum h is classified to the l^{th} income class and $y_{khl} = 0$ otherwise and m_h the total number of households selected in stratum h . In the RIS $L = 10$. Income distributions for persons are defined and estimated similarly to (4.1), (4.2), with M_{lr} the number of persons from area r , belonging to the l^{th} income category. The HT estimator for M_{lr} is now defined as

$$\hat{M}_{lr} = \sum_{h \in r} \sum_{k=1}^{m_h} \frac{1}{\pi_k} \sum_{j=1}^{N_k} y_{kjhl},$$

where $y_{kjhl} = 1$ if person j from household k and stratum h is classified to the l^{th} income class and $y_{kjhl} = 0$ otherwise.

For sample size determination, precision specifications for the estimated income distributions are required. For stratified sampling designs, Neyman allocations are often considered to determine minimum sample sizes and optimal allocations to meet precision requirements at aggregated levels (Cochran 1977). Power allocations are useful to find the right balance between precision requirements for aggregates and strata (Bankier 1988). In this application the minimum sample size is based on precision requirements for the individual strata, i.e., neighbourhoods, which is the most detailed publication level.

If precision requirements are specified for the separate classes of the income distributions, then the income class with the largest population variance determines the minimum required sample size, resulting in unnecessarily large sample sizes. As an alternative the square root of the mean over the variances of the estimated income classes of an income distribution is proposed as a precision measure for the estimated income distributions. With this measure the influence of the most imprecise income class on the minimum sample size will be reduced. The square root of the mean over the variances of the estimated income classes of an income distribution is called the average standard error measure and is defined as

$$s = \sqrt{\frac{1}{L} \sum_{l=1}^L V(\hat{P}_{lr})}. \quad (4.3)$$

In this section an exact expression for s will be derived as well as an approximation that can be used to estimate the minimum required sample size which does not require information about income distributions or variances.

Since neighbourhoods are the most detailed areas for which income distributions are published, precision requirements for sample size determination are specified at this level. Since neighbourhoods are used as the stratification variable in the sample design, expressions for s can be derived under simple random sampling without replacement of core persons within each neighbourhood. It is proved in the appendix that an expression for the average standard error measure s_h in (4.3) for an income distribution is given by

$$s_h = \sqrt{\frac{1}{L} \frac{N_h - n_h}{n_h} \frac{1}{N_h - 1} \left(\frac{N_h}{M_h^2} \sum_{l=1}^L \sum_{k=1}^{M_h} \frac{y_{khl}}{g_{kh}} - \sum_{l=1}^L \left(\frac{M_{lh}}{M_h} \right)^2 \right)}, \quad (4.4)$$

with M_h the number of households in stratum h and M_{lh} the number of households in stratum h belonging to the l^{th} income class. Note that if $g_{kh} = 1$ for all households in the population of stratum h , then it follows that $M_h = N_h$ and that formula (4.1) simplifies to

$$V(\hat{P}_h) = \frac{N_h - n_h}{n_h} \frac{1}{N_h - 1} (P_h (1 - P_h)),$$

which can be recognized as the variance of an estimated fraction under simple random sampling without replacement (Cochran 1977, Chapter 3).

Minimum sample size requirements based on (4.4) require information about the income distribution and its variance from preceding periods. Since this information is generally not available at the design phase of a panel, it is useful to have an upper bound for the average standard error measure for the income distribution in (4.4). This is comparable to taking the variance for a parameter defined as a proportion, which reaches a maximum when the proportion is 0.5 for calculating the minimum sample size for a survey. It is shown in the appendix that an upper bound for the average standard error measure s_h for an income distribution, specified in (4.4) is given by

$$s_h \leq \sqrt{\frac{1}{L} \frac{N_h - n_h}{n_h} \frac{1}{N_h - 1} \left(\frac{N_h}{M_h^2} \sum_{t=1}^T \frac{M_{th}}{t} - \frac{1}{L} \right)}, \quad (4.5)$$

with M_{th} the number of households of size t in stratum h .

If $g_{kh} = 1$ for all households in the population of stratum h and the number of classes of the income distribution $L = 2$, then it follows that the approximation for the average standard error measure s_h in (4.5) can be simplified to

$$s_h \leq \sqrt{\frac{N_h - n_h}{n_h} \frac{1}{(N_h - 1)} \frac{1}{4}},$$

which equals the square root of the maximum variance of an estimated fraction at $\hat{P} = 0.5$ under simple random sampling. This illustrates that the approximation for the average standard error measure in (4.5) can be interpreted as a generalization of the approximation of the maximum variance of an estimated fraction at $\hat{P} = 0.5$, often used in sample size determination. The average standard error measure has its maximum value in the case of an equal distribution of the households over the income categories, i.e., $\hat{P}_h = 1/L$ for $l = 1, \dots, L$. In this situation the approximation for s_h is exact, which follows directly from equation (4.3).

Equating the expression for s_h in (4.5) to a pre-specified maximum value, say Δ_h , results in the following expression for the minimum sample size of core persons

$$n_h \geq \frac{\left(\frac{N_h}{M_h} \right)^2 \sum_{t=1}^T \frac{M_{th}}{t} - \frac{N_h}{L}}{(N_h - 1) L \Delta_h^2 + \frac{N_h}{M_h^2} \sum_{t=1}^T \frac{M_{th}}{t} - \frac{1}{L}}. \quad (4.6)$$

The information required to estimate the minimum sample size is the total number of persons and the total number of equally sized households for neighbourhoods. No information about the expected income distribution or its variance is required. More precise estimates for the minimum sample size can be obtained with the expression in (4.4), but require sample information from, for example, previous periods about the income distributions.

Expression (4.6) gives the minimum sample size for core persons. Subsequently all household members of each core person are included in the sample. As a result, households can be included in the sample more than once and the sample size in terms of unique households and unique persons is random. To plan a survey and control survey costs, it is necessary to know the expected number of unique households and unique persons if a sample of core persons of size n_h is drawn. In the appendix it is proved that the expected number of unique households in a sample of n_h core persons, drawn by means of simple random sampling without replacement from a finite population of size N_h is given by

$$D_h = \sum_{t=1}^T M_{th} \left(1 - \frac{\prod_{i=0}^{t-1} (N_h - n_h - i)}{\prod_{i=0}^{t-1} (N_h - i)} \right). \quad (4.7)$$

The expected number of unique persons in a sample of n_h core persons, drawn by means of simple random sampling without replacement from a finite population of size N_h follows directly from (4.7) and is given by

$$D_h^{[p]} = \sum_{t=1}^T tM_{th} \left(1 - \frac{\prod_{i=0}^{t-1} (N_h - n_h - i)}{\prod_{i=0}^{t-1} (N_h - i)} \right). \quad (4.8)$$

Since the expected numbers of unique households and persons are random variables, it would be useful to have an uncertainty measure for these expected values. Variance expressions for (4.7) and (4.8) are however not straightforward and therefore left for further research.

Sample size calculations are conducted at the level of neighbourhoods. It was finally decided to select core persons with a sampling fraction of 0.16. With this sample size, the maximum value for the average standard error measure s_h at the level of neighbourhoods amounts to about 0.01 for the estimated household income distributions. With a total population of about 12 million persons, this resulted in a sample size of about 2.1 million core persons and an expected sample size of about 4.6 million unique persons. This sample was drawn in 1994, which was the start of the panel for the Dutch RIS.

5 Linear weighting

For household surveys like the RIS, estimates are required for person characteristics as well as household characteristics. Let t_y denote the total of a target variable y . With linear weighting, an estimator for a person based target variable is defined as

$$\hat{t}_y = \sum_{h=1}^H \sum_{k \in 1}^{m_h} \sum_{j \in k} w_{kj} y_{kjh}, \quad (5.1)$$

with y_{kjh} the value of the target variable for person (k, j, h) and w_{kj} a weight for person j belonging to household k . An estimator for a household based target variable is given by

$$\hat{t}_y = \sum_{h=1}^H \sum_{k=1}^{m_h} w_k y_{kh}, \quad (5.2)$$

with y_{kh} the value of the target variable for household k from stratum h and w_k a weight for the corresponding household.

Weights are obtained by means of the GREG estimator to use auxiliary variables which are observed in the sample and for which the population totals are known from other sources (Särndal et al. 1992). Consequently, the weights reflect the (unequal) inclusion expectations of the sampling units and an adjustment such that for auxiliary variables the weighted observations sum to the known population totals. Often categorical variables like gender, age, marital status or region are used as auxiliary variables. Due to the fact that the values of auxiliary variables differ from person to person within the same household, different weights can be derived for people from the same household. To ensure that relationships between household variables and person variables are reflected in estimated totals, it is relevant to apply a weighting method which yields one unique household weight for all its household members. If the weights for persons within a household are the same, then household and person based estimates of the same target variables are consistent with each other (for example the total income estimated from households and that from persons). This can be achieved with so-called integrated weighting methods.

Lemaître and Dufour (1987) apply an integrated weighting method at the persons level and replace the original auxiliary variables defined at the person level by the corresponding household mean. In this way, members of the same household have the same inclusion expectation and share the same auxiliary information, and therefore the resulting regression weights are forced to be the same. Nieuwenbroek (1993) proposes a slightly more general approach by applying the linear weighting method at the household level, where the auxiliary information of person based characteristics is aggregated at the household level. Nieuwenbroek (1993) mentions that the linear weighting method at the household level is equal to the linear weighting method of Lemaître and Dufour (1987) at the person level, if the residual variance of the regression model at the household level is chosen proportional to the number of persons within the household. Steel and Clark (2007) and Estevao and Särndal (2006) further generalize the integrated weighting of person and household surveys. Steel and Clark (2007) address the issue of whether the cosmetic benefits of integrated weighting result in an increased design variance of the GREG estimates. They show that large-sample design variances obtained by linear weighting at the household level is less than or equal to the design variance obtained with linear weighting at the person level. For small samples there can be a small increase in the design variance due to integrated weighting. As a result there is little or no loss in efficiency by applying an integrated weighting method.

In this paper the integrated weighting approach at the household level is applied. Let \mathbf{x}_{kh} denote a q -vector containing q auxiliary variables for household k from stratum h . Person based characteristics are aggregated to household totals. The GREG estimator is derived from a linear regression model that specifies the relation between the target variable and the available auxiliary variables for which population totals are known, and is defined as:

$$y_{kh} = \mathbf{x}_{kh}^t \boldsymbol{\beta} + e_{kh}, \quad \text{with} \quad E_m(e_{kh}) = 0, \quad V_m(e_{kh}) = \sigma_{kh}^2. \quad (5.3)$$

In (5.3) $\boldsymbol{\beta}$ denotes a vector containing the q regression coefficients of the regression of y_{kh} on \mathbf{x}_{kh} and e_{kh} the residuals and E_m and V_m denote the expectation and variance with respect to the regression model. In this application, the variance structure is taken proportional to the household size, i.e., $\sigma_{hk}^2 = g_k \sigma^2$. Nieuwenbroek (1993) shows that in this case the weighting applied at the household level is equal to the method of Lemaître and Dufour (1987).

Regression weights for the households are finally obtained by

$$w_k = \frac{1}{\pi_k} \left(1 + (\mathbf{t}_x - \hat{\mathbf{t}}_{x\pi})^t \left(\sum_{k=1}^m \frac{\mathbf{x}_{kh} \mathbf{x}_{kh}^t}{\pi_k g_k} \right)^{-1} \frac{\mathbf{x}_{kh}}{g_k} \right),$$

with \mathbf{t}_x a q vector containing the known population totals of the auxiliary variables \mathbf{x} , $\hat{\mathbf{t}}_{x\pi}$ the HT estimator for \mathbf{t}_x . The weights calculated at the household level can be used for weighting person based characteristics of the corresponding household members, using formula (5.1) since $w_{kj} = w_k$ for all persons belonging to the same household k .

6 Variance estimation

Parameters of the RIS are estimated as the ratio of two population totals

$$\hat{R} = \frac{\hat{t}_y}{\hat{t}_z}, \quad (6.1)$$

where \hat{t}_y and \hat{t}_z are GREG estimators defined by (5.1) or (5.2) in the case of person-based or household-based target variables, respectively. The variance of (6.1) under a sample design where core persons are drawn by means of stratified simple random sampling, and all household members of these core persons are included in the sample can be approximated by

$$V(\hat{R}) = \frac{1}{\hat{t}_z^2} \sum_{h=1}^H \frac{N_h^2 (1 - f_h)}{n_h} \frac{1}{N_h - 1} \sum_{k=1}^{N_h} \left(\frac{e_{kh}}{g_k} - \frac{1}{N_h} \sum_{k'=1}^{N_h} \frac{e_{k'h}}{g_{k'}} \right)^2, \quad (6.2)$$

where $f_h = n_h / N_h$, $e_{kh} = (y_{kh} - \mathbf{x}_{kh}^t \mathbf{b}_y) - R(z_{kh} - \mathbf{x}_{kh}^t \mathbf{b}_z)$, and \mathbf{b}_y and \mathbf{b}_z are the finite population regression coefficients of the regression of y_{kh} and z_{kh} respectively on \mathbf{x}_{kh} . An estimator for the variance specified in (6.2) is given by

$$\hat{V}(\hat{R}) = \frac{1}{\hat{t}_z^2} \sum_{h=1}^H (1 - f_h) \frac{n_h}{n_h - 1} \sum_{k=1}^{n_h} \left(w_k \hat{e}_k - \frac{1}{n_h} \sum_{k'=1}^{n_h} w_{k'} \hat{e}_{k'h} \right)^2, \quad (6.3)$$

where $\hat{e}_{kh} = (y_{kh} - \mathbf{x}_{kh}^t \hat{\mathbf{b}}_y) - \hat{R}(z_{kh} - \mathbf{x}_{kh}^t \hat{\mathbf{b}}_z)$ and $\hat{\mathbf{b}}_y$ and $\hat{\mathbf{b}}_z$ are the HT type estimators for \mathbf{b}_y and \mathbf{b}_z . These results follow directly from inserting first and second order inclusion expectations specified in (3.3) through (3.6) in the general approximation for the variance of the ratio of two GREG estimators and its estimator (Särndal et al. 1992, Section 7.13).

The same expressions for the variance can be derived from the variance expressions proposed for the Generalized Weight Share method in the case of indirect sampling. In Lavallée (1995), variance expressions for the HT estimator are based on the sampling design used to select the sample s^A of n units from population U^A with transformed target variables, say z_i . In this application each unit in U^A has exactly one link with a unit in U^B . As a result z_i in Lavallée (1995) is in this case defined as the sum over the target variables of all elements in cluster k , divided by the number of units in cluster k with a link to population U^A , i.e., $z_i = y_k / g_k$ for all $i \in U^A$ that have a link with cluster $k \in U^B$. Inserting the first and second order inclusion probabilities for stratified simple random sampling without replacement and the transformed variables z_i (where the target variable y_k is replaced by the residual of the regression on the cluster totals e_k) in the variance formula for a ratio gives (6.2). Result (6.3) follows in a similar way.

7 Application

In the RIS, core persons are selected from the population aged 15 years and older through stratified simple random sampling without replacement with a sample fraction of 0.16. In this application results are presented for a large municipality (Rotterdam), a municipality of intermediate size (Enschede) and a small municipality (Sevenum) for three consecutive years 2006, 2007 and 2008. Population and sample sizes for these three municipalities are summarized in Table 7.1.

Table 7.1
Population and sample size RIS for three Dutch municipalities

Municipality	Population		Sample		
	Households	Persons 15 and older	Core persons	Unique households	Unique persons
Rotterdam	293,400	484,000	73,000	67,600	171,400
Enschede	74,200	128,000	19,300	17,600	46,300
Sevenum	2,950	6,100	870	750	2,500

Target variables of interest for the RIS are:

- Income distribution of households in ten classes where the categories are based on ten percentage point quantiles (deciles) of the national distribution using standardized household income (abbreviated as IncDistHh);
- Mean standardized household income (abbreviated as HHinc);
- Mean disposable income of persons that receive income during the 52 weeks of the year (abbreviated as Pinc).

Disposable income of a person is total income of a person minus his or her current taxes. Total income contains earnings, profit, income from capital and savings, and social or other benefits. Standardized household income is defined as the total disposable income of a household corrected for differences in household size and composition. In the literature, this is also known as the equalised spendable income (OECD 2013).

Estimates for official publications of the RIS are obtained with the GREG estimator using the method of Lemaître and Dufour (1987). Since this survey does not suffer from nonresponse, auxiliary information is used in the estimation for variance reduction and consistency between the marginals of different publication tables. Inclusion expectations are based on the formulas derived in Subsection 3.1. For each municipality the following weighting scheme is applied in the GREG estimator:

$$\text{Age (7)} \times \text{Gender} + \text{Age (4)} \times \text{Gender} \times \text{MaritalStatus (2)} + \text{Address (2)} \times \text{HHsize (5)}.$$

All auxiliary variables are categorical. The numbers between brackets denote the number of categories. MaritalStatus distinguishes between people who are married and other forms of marital status. Address distinguishes between addresses where one family is residing and other types of addresses. HHsize stands for household size and distinguishes between households with one, two, three, four, and five or more persons. Estimates for HHinc and Pinc with their standard errors based on the HT estimator, the GREG estimator and the GREG estimator with the method of Lemaître and Dufour (1987) are given in Table 7.2. In Figure 7.1 the income distributions IncDistHh estimated with the HT estimator, GREG estimator and the GREG estimator with the method of Lemaître and Dufour (1987) are plotted with a 95% confidence interval for Rotterdam and Sevenum in 2008. The standard errors for these estimates are compared in a separate histogram. In Figure 7.2 the IncDistHh for Rotterdam and Sevenum estimated with the method of Lemaître and Dufour (1987) are given for 2006, 2007 and 2008. See van den Brakel (2013) for more detailed output of the income distributions.

Table 7.2
Estimation results RIS for Rotterdam (large city), Enschede (intermediate city), and Sevenum (small village), standard errors in brackets

	Variable	Year	HT		GREG		GREG consistent (L&D)	
Rotterdam	HHinc	2006	19,790	(83)	20,134	(80)	20,161	(76)
		2007	22,306	(73)	22,950	(64)	22,866	(64)
		2008	23,750	(78)	24,511	(69)	24,410	(68)
	Pinc	2006	22,074	(94)	22,219	(84)	22,233	(93)
		2007	24,094	(82)	24,362	(75)	24,432	(78)
		2008	25,325	(84)	25,625	(75)	25,705	(78)
Enschede	HHinc	2006	19,810	(128)	20,353	(111)	20,300	(107)
		2007	20,878	(128)	21,716	(107)	21,753	(105)
		2008	22,254	(148)	23,235	(125)	23,237	(123)
	Pinc	2006	20,402	(102)	20,608	(92)	20,590	(92)
		2007	21,387	(115)	21,751	(103)	21,852	(106)
		2008	22,235	(123)	22,659	(110)	22,724	(114)
Sevenum	HHinc	2006	25,696	(799)	25,698	(734)	25,968	(711)
		2007	28,207	(618)	28,901	(520)	29,026	(490)
		2008	31,466	(795)	32,372	(715)	32,536	(694)
	Pinc	2006	21,328	(466)	21,680	(428)	21,712	(428)
		2007	24,056	(456)	24,219	(396)	24,459	(393)
		2008	24,980	(468)	25,482	(426)	25,644	(455)

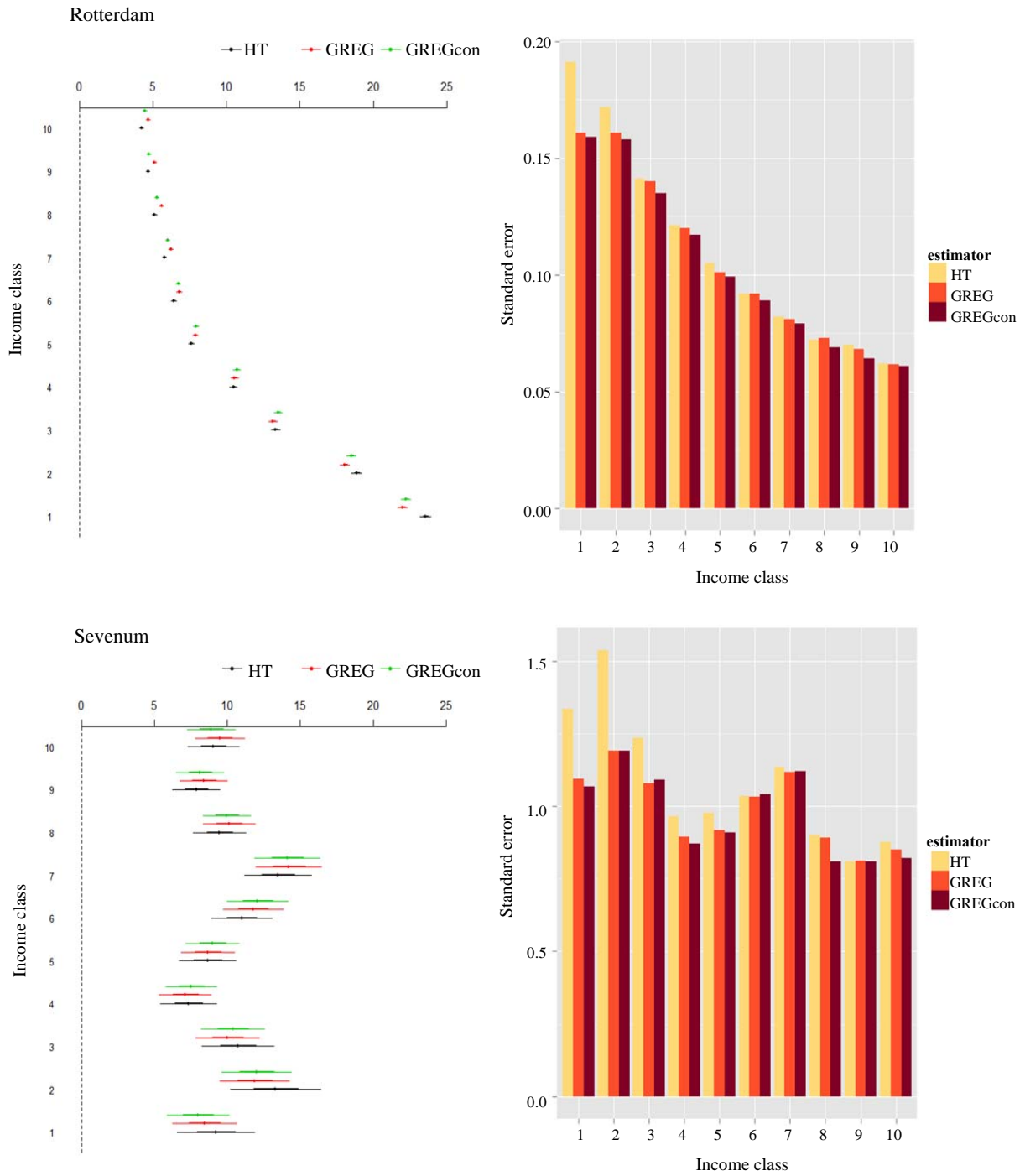


Figure 7.1 IncDistHh in percentages for Rotterdam and Sevenum (left panels) with Horvitz-Thompson estimator, GREG estimator and integrated GREG estimator (GREGcon), with 95% confidence intervals. Standard errors of the corresponding estimators are plotted in the right panels.

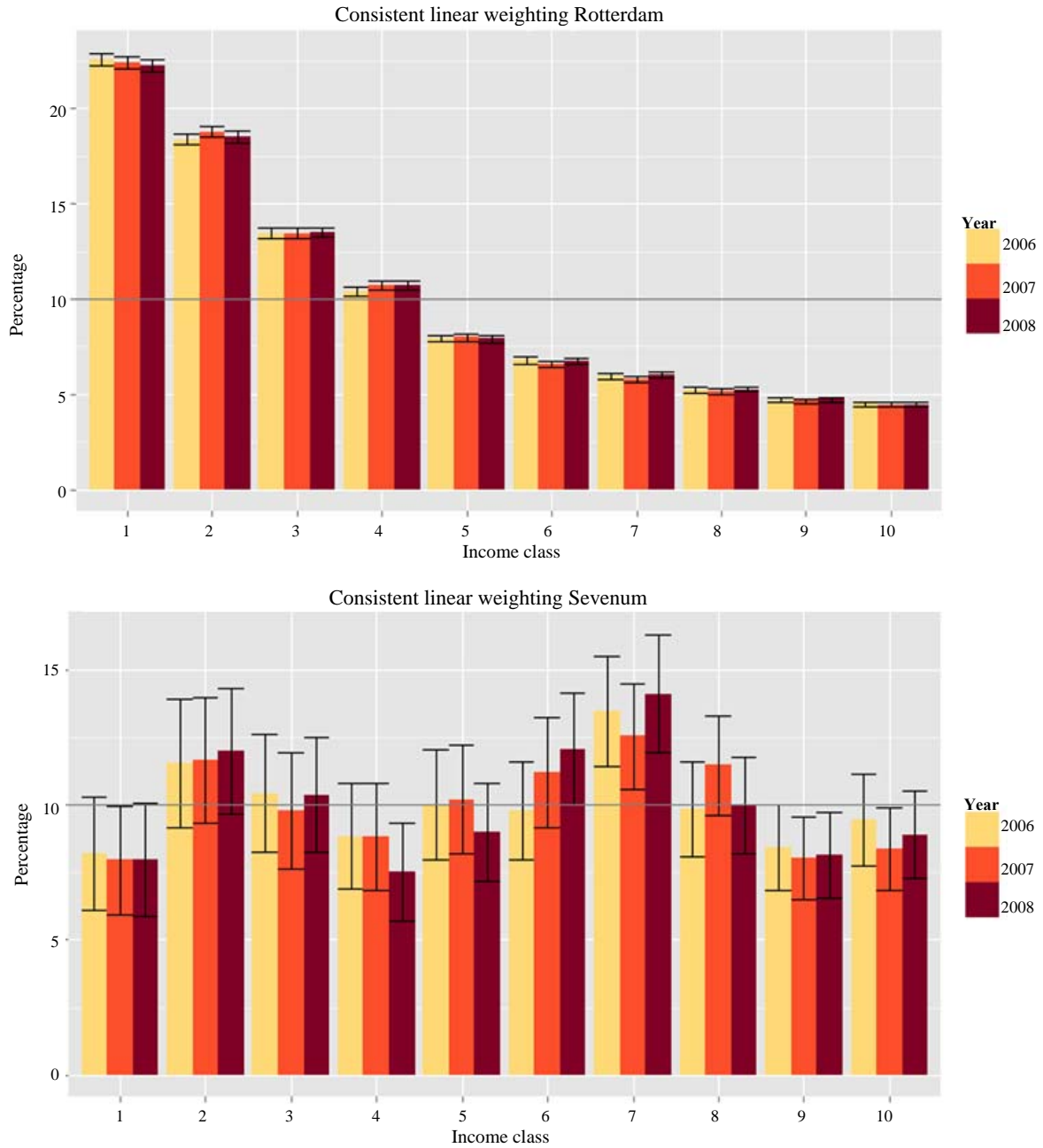


Figure 7.2 IncDistHh in percentages for Rotterdam (upper panel) and Sevenum (lower panel) estimated with integrated weighting for 2006, 2007 and 2008 with 95% confidence intervals. Grey line refers to the national income distribution.

The observed income distributions in Figures 7.1 and 7.2 are a result of the demographic compositions in both municipalities. Rotterdam is a city where the fraction of households in low income categories are above the national average, since the fractions in the first three categories are above 10%. The fraction of households in higher income categories, on the other hand, are below the national average, since these fractions are below 10%. This is a typical distribution for a large university city with a high fraction of

non-western immigrants. Sevenum on the other hand is a small village close to a large industrial city. Such villages typically have small fractions of immigrants, no students and large fractions of households with one or two people that receive income during 52 weeks of the year. This explains why the fraction of households in the lowest income category is below the national average and the fraction of households in the higher income categories (6, 7 and 8) is above the national average. Sevenum is a village that does not attract extreme rich households.

Since HHinc and Pinc are based on different income definitions and since Pinc is the average over the domains of people that receive income during 52 weeks of the year, the differences between the two means vary between municipalities. For a large university city like Rotterdam, the mean standardized household income is typically smaller compared to the mean of disposable personal income averaged over people that receive income during 52 weeks of the year. Other cities with large universities show a similar picture. In a small but rich village like Sevenum, the situation is the other way around.

Another remarkable result is that in Rotterdam and Enschede the difference between the HT estimator and the GREG estimator is relatively large compared to the standard errors. Given the large sample size and the fact that there is no nonresponse, these differences are expected to be smaller. A possible explanation is that Rotterdam and Enschede are large university cities. Students are often identified in the tax register (used as the sample frame) in a different way than they appear in the population register (used to derive population distributions of the auxiliary variables), in particular with respect to their household situation.

For each municipality there is a steady increase over time in the mean of the income for households and persons. Also the income distributions for each municipality show a stable pattern over the years. This can be expected if a panel is applied in combination with large sample sizes to estimate phenomena that are not very volatile in time.

Comparing GREG estimates with and without using the method of Lemaître and Dufour (1987) shows that standard errors of estimated household parameters are smaller if the method of Lemaître and Dufour (1987) is applied. This is particularly visible for the mean household income in the small sample of Sevenum. For estimated person based parameters, on the other hand, the method of Lemaître and Dufour (1987) slightly increases the standard error compared to the regular GREG estimator. This suggests that the assumed variance structure for the residuals in the underlying regression model in the case of integrated weighting better fits the household-based variables than the person-based variables.

8 Discussion

Households, due to their instability over time, are inappropriate as sampling units in panels conducted to collect information at the level of households or persons. In this paper, a sample design is proposed where persons are drawn through a self-weighted sample design. At each point in time, the household members of these so-called core persons are included in the sample. This results in a sample where households can be drawn more than once but with a maximum that is equal to the household size. Households are included with expectations proportional to the household size. First and second order inclusion expectations for households are derived under an equal probability sample design for selecting core persons. These inclusion expectations can be used in a similar way to the more common inclusion probabilities in design-based and model-assisted inference.

The sample design in this paper is a special case of indirect sampling (Lavallée 1995, 2007). In the case of a self-weighted sample design it is shown that first and second order inclusion expectations for this sample design can be derived in a relatively straightforward manner from the household composition of the core persons at each point in time. In the case of more complex sample designs, the Generalized Weight Share method (Lavallée 1995, 2007), is required to construct inclusion weights at each point in time.

The advantage of the proposed sample design is that the estimation procedure is simpler than the Generalized Weight Share method. The design is particularly useful if core persons are selected with a self-weighted sampling design. If, due to, e.g., minimum precision and maximum cost requirements, an unequal probability design for the selection of core persons is required, then the Generalized Weight Share method is required. Since core persons remain in the panel indefinitely, this sample design is particularly appropriate for register-based household panels where all the required information is derived from administrative data. For interview-based household panels some kind of rotating design is required to cope with problems like panel attrition.

In the paper the so called average standard error measure, defined as the square root of the mean over the variances of the estimated income classes of an income distribution, is proposed as a precision measure for minimum sample size determination. It is shown that the maximum value of this precision measure corresponds with a distribution where the proportions in the categories are equal. It is also shown that this result can be seen as generalization of the variance of a fraction taking its maximum value at 0.5. An expression for the minimum required sample size to meet a pre-specified precision for estimated distributions is derived. Since households can be included more than once in the sample, an expression for the expected number of unique households in a sample is also derived.

A topic for further research is to combine this mean standard error measure with a Neyman allocation or power allocations to have expressions for the minimum sample size based on precision requirements for estimated distributions at aggregates of strata. This results in an unequal inclusion probability design for the core persons and requires the Generalized Weight Share method for deriving appropriate weights.

In the context of household surveys and panels, weighting procedures that enforce equal regression weights for persons within the same household are relevant in order to enforce consistency between person based and household based estimates. In this paper an integrated weighting approach based on Lemaître and Dufour (1987) is applied to the RIS. In this application standard errors obtained with Lemaître and Dufour (1987) are smaller than a non-integrated weighting procedure for household based estimates. For person based estimates, standard errors can be slightly larger. These results are in line with Steel and Clark (2007), who show that the large-sample design variance of integrated weighting at the household level is smaller than or equal to the design variance obtained with non-integrated weighting at the person level. In their simulation they also report small increases of the design-variances due to integrated weighting in the case of small sample sizes.

Integrated weighting of Lemaître and Dufour (1987) at the household level is obtained by assuming a variance structure for the residuals that is proportional to the household size (Nieuwenbroek 1993). If household characteristics are proportional to household size, then it can be anticipated that such a variance structure better explains the variation of the household variables in the population compared to a variance structure that assumes equal residual variance for the households. For person based variables such a variance

structure might be less efficient but the additional advantage of integrated weighting is that totals for household and person based income, which can be derived directly from their means, are consistent.

Acknowledgements

The views expressed in this paper are those of the author and do not reflect the policies of Statistics Netherlands. The author is grateful to the Associate Editor and the unknown referees for giving constructive comments on two former drafts of the paper. The author also thanks Drs. M. van den Brakel-Hofmans for making the RIS data available.

Technical appendix

Proof of equation (4.4)

An expression for the variance of the estimated fraction of households in income class l can be derived from the general expression for the variance of the HT estimator (Särndal et al. 1992, Section 2.8):

$$V(\hat{P}_{th}) = \frac{1}{M_h^2} \sum_{k=1}^{M_h} \sum_{k'=1}^{M_h} (\pi_{kk'h} - \pi_{kh} \pi_{k'h}) \frac{y_{khl}}{\pi_{kh}} \frac{y_{k'hl}}{\pi_{k'h}}. \quad (\text{A.1})$$

Inserting first and second order inclusion expectations specified in (3.3) through (3.6), and taking advantage of the property that $y_{khl} = y_{khl}^2$ since the values of the target variable are restricted to zero or one, it follows after some algebra that (A.1) can be simplified to

$$V(\hat{P}_{th}) = \frac{N_h - n_h}{n_h} \frac{1}{N_h - 1} \left(\frac{N_h}{M_h^2} \sum_{k=1}^{M_h} \frac{y_{khl}}{g_{kh}} - \left(\frac{M_{lh}}{M_h} \right)^2 \right). \quad (\text{A.2})$$

Result (4.4) is obtained by inserting (A.2) into (4.3).

Proof of equation (4.5)

The *population* of households in stratum h can be divided into T subpopulations of equally sized households. Let M_{th} denote the number of households of size t in stratum h . Now it follows for the double summation between brackets for the expression of s in (4.4) that

$$\sum_{l=1}^L \sum_{k=1}^{M_h} \frac{y_{khl}}{g_{kh}} = \sum_{l=1}^L \sum_{t=1}^T \sum_{k=1}^{M_{th}} \frac{y_{khl}}{t} = \sum_{t=1}^T \frac{M_{th}}{t}. \quad (\text{A.3})$$

According to the Cauchy-Schwartz inequality (Cochran 1977, Section 5.5) it follows for the single summation between brackets for the expression of s_h in (4.4) that

$$\sum_{l=1}^L \left(\frac{M_{lh}}{M_h} \right)^2 = \sum_{l=1}^L P_{lh}^2 \geq \frac{1}{L}. \quad (\text{A.4})$$

Result (4.5) is obtained by inserting (A.3) and (A.4) in the expression for s in (4.4).

Proof of equation (4.7)

Let $\tilde{\pi}_{tkh}$ denote the inclusion probability for household k from stratum h of size t . Since equally sized households share the same first order probabilities, it follows that $\tilde{\pi}_{ikh} = \tilde{\pi}_{ik'h} \equiv \tilde{\pi}_{th}$. Let I_{tkh} denote an indicator variable, taking value 1 if household k from stratum h of size t is included in the sample and zero otherwise. The expected number of unique households can be derived as

$$\begin{aligned} D_h &= E\left(\sum_{t=1}^T \sum_{k=1}^{M_{th}} I_{tkh}\right) = \sum_{t=1}^T M_{th} \tilde{\pi}_{th} \\ &= \sum_{t=1}^T M_{th} \left(1 - \frac{\binom{N_h - t}{n_h}}{\binom{N_h}{n_h}}\right) = \sum_{t=1}^T M_{th} \left(1 - \frac{(N_h - n_h)(N_h - n_h - 1) \dots (N_h - n_h - t + 1)}{N_h (N_h - 1) \dots (N_h - t + 1)}\right). \end{aligned}$$

References

- Bankier, M.D. (1988). Power allocations: Determining sample sizes for subnational areas. *The American Statistician*, 42, 174-177.
- Bethlehem, J.G. (2009). *Applied Survey Methods*, New Jersey: John Wiley & Sons, Inc.
- Cochran, W.G. (1977). *Sampling Techniques*, New York: John Wiley & Sons, Inc.
- Deville, J.-C., and Lavallée, P. (2006). Indirect sampling: The foundations of the generalized weight share method. *Survey Methodology*, 32, 2, 165-176.
- Ernst, L. (1989). Weighting issues for longitudinal household and family estimates. In *Panel Surveys*, (Eds. D. Kasprzyk, G. Duncan, G. Kalton and M.P. Singh). New York: John Wiley & Sons, Inc., 135-159.
- Estevao, V.M., and Särndal, C.-E. (2006). Survey estimates by calibration on complex auxiliary information. *International Statistical Review*, 74, 127-147.
- Horvitz, D.G., and Thompson, D.J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47, 663-685.
- Kalton, G., and Brick, J.M. (1995). Weighting schemes for household panel surveys. *Survey Methodology*, 21, 1, 33-44.
- Lavallée, P. (1995). Cross-sectional weighting of longitudinal surveys of individuals and households using the weight share method. *Survey Methodology*, 21, 1, 25-32.
- Lavallée, P. (2007). *Indirect Sampling*, New York: Springer Verlag.
- Lemaître, G., and Dufour, J. (1987). An integrated method for weighting persons and families. *Survey Methodology*, 13, 2, 199-207.

- Lynn, P. (2009). Methods for longitudinal surveys. In *Methodology of Longitudinal Surveys*, (Ed., P. Lynn), Wiley, Chichester, 1-19.
- Narain, R. (1951). On sampling without replacement with varying probabilities. *Journal of the Indian Society of Agricultural Statistics*, 3, 169-174.
- Nieuwenbroek, N.J. (1993). An integrated method for weighting characteristics of persons and households using the linear regression estimator. Research paper, BPA nr: 8555-93-M1-1, Statistics Netherlands, Heerlen.
- OECD (2013). *OECD Framework for Statistics on the Distribution of Household Income, Consumption and Wealth*. OECD publishing, <http://dx.doi.org/10.1787/9789264194830-en>.
- Särndal, C.-E., Swensson, B. and Wretman, J. (1992). *Model Assisted Survey Sampling*. New-York: Springer-Verlag.
- Smith, P., Lynn, P. and Elliot, D. (2009). Sample design for longitudinal surveys. In *Methodology of Longitudinal Surveys*, (Ed., P. Lynn), Wiley, Chichester, 21-33.
- Steel, D.G., and Clark, R.G. (2007). Person-level and household-level regression estimation in household surveys. *Survey Methodology*, 33, 1, 51-60.
- van den Brakel, J.A. (2013). Sampling and estimation techniques for household panels. Discussion paper 2013-15, Statistics Netherlands, Heerlen. <http://www.cbs.nl/NR/rdonlyres/B4F85FB9-52F2-4B8A-94C4-56DA43F2250D/0/201315x10pub.pdf>.
- Wallgren, A., and Wallgren, B. (2007). *Register-Based Statistics: Administrative Data for Statistical Purposes*. New York: John Wiley & Sons, Inc.