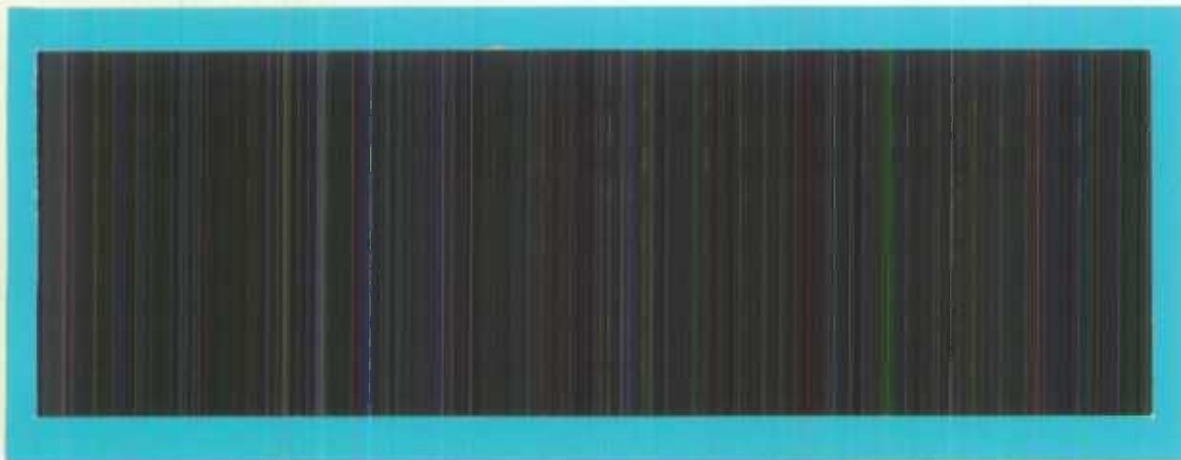




Statistics
Canada

Statistique
Canada



Methodology Branch

Social Survey
Methods Division

Direction de la méthodologie

Division des méthodes
d'enquêtes sociales

11-613

no. 91-02

11-613

C. 2

11-613

Canada



BOYHK

C.2

WORKING PAPER NO. SSMD-91-002 E

METHODOLOGY BRANCH



**EFFECTS OF SAMPLE OVERLAP ON THE
ESTIMATION AND DETECTION OF CHANGE IN INCOME**

SSMD-91-002 E

**Edward J. Chen and S. Kumar
Social Survey Methods Division
Statistics Canada**

January, 1991

ABSTRACT

It is well known that an improved estimate of change, from one time point to the next, can be obtained if positive correlation between cross-sectional estimates can be induced by sample overlap between those time points. However, in the case of survey with multistage sample designs, such as Survey of Consumer Finances, the gains in efficiencies (due to sample overlap) will depend on the stage at which the sample overlap occurs and the proportion of overlap. We examine the effects of sample overlap on the estimation and detection of change in income by using the data from Survey of Consumer Finances conducted in 1988 and 1989.

Key Words: multistage sample design, overlap sample, relative efficiency, correlation, detection of change.

RÉSUMÉ

Il est bien connu que l'estimation du changement entre deux points dans le temps peut être améliorée en provoquant, par le chevauchement des échantillons entre ces points, une corrélation positive entre des estimations transversales. Toutefois, dans le cas des enquêtes ayant des plans de sondage à plusieurs degrés, comme l'enquête sur les finances des consommateurs, les gains d'efficacité (dus au chevauchement des échantillons) dépendront du degré auquel s'effectue le chevauchement et de la proportion du chevauchement. Nous examinons les effets du chevauchement des échantillons sur l'estimation et la détection du changement dans les revenus à partir des données de l'enquête sur les finances des consommateurs réalisée en 1988 et en 1989.

Mots clés: Plan de sondage à plusieurs degrés, échantillons chevauchants, efficacité relative, corrélation, détection du changement.

TABLE OF CONTENTS

	Page
1. Introduction	1
2. Notation	4
3. Relative Efficiencies of Two Designs	6
3.1 Average Income of E.F. 2+	6
3.2 Average Income of Unattached Individuals	7
4. Detection of Change in Average Incomes	8
5. Conclusions	11
Table I Correlation Coefficients and Relative Efficiencies for E.F. 2+	14
Table II Correlation Coefficients and Relative Efficiencies for Unattached Individuals	15
Table III Table of 't-values' and 'P-Values' for E.F. 2+	16
Table IV Table of 't-values' and 'P-Values' for Unattached Individuals	17
Appendix I - Estimation of Correlation Coefficients	18
Appendix II - Estimation of the Change and Variance	20
Acknowledgements	22
References	22

1. INTRODUCTION

Statistics Canada has been conducting the Survey of Consumer Finances (SCF) annually since 1971. This survey involves the collection and analysis of data about individual, family, and household incomes and the sources of these incomes. The SCF is a supplemental survey to the monthly Labour Force Survey (LFS). Thus the target population for SCF is the same as for the LFS. The sample represents virtually all private households in Canada except for households located on Indian reserves and inmates of institutions. The residents of Yukon and Northwest Territories are also excluded from the survey.

The LFS sample consists of six panels or rotation groups (See [3] for the methodology). Under the rotation scheme, each month one-sixth of the households rotate out of the sample and are replaced by new households. SCF samples normally use four of the six LFS rotation groups. However, during each of the 1988 and 1989 surveys, five rotation groups were used.

It is well known that improved estimates of change, from time 1 to time 2, can be obtained if it is based on cross-sectional estimates that are positively correlated (as opposed to those based on non-correlated cross-sectional estimates). For recurring surveys, this positive correlation between two cross-sectional estimates is induced by having some common units in the two samples for points 1 and 2. In multistage stratified sample surveys, such as LFS and SCF, overlap of sampled units may occur at the ultimate stage of sampling or at a prior stage but not at the ultimate stage. In the case of LFS:

(i) matching of some household units (the ultimate sampled units) at two points in time implies corresponding matching at the primary sampling unit (PSU)

level (a sampled unit at an earlier stage of sampling) at those two times. This is the case when the two points are at most five months apart; and

(ii) matching of sampled PSU's does not necessarily imply common sampled households within those PSU's. This situation occurs when the two points are at least six months apart.

The magnitude of the induced correlation between cross-sectional estimates will depend upon the matching level of the selected units, i.e. whether it is at the household level or at the PSU level (without matching households). Henceforth, the term 'PSU overlap' will mean that there are some identically selected PSU's but no matching (overlap) of sampled households within those PSU's.

Intuitively, it is expected that the correlation between two cross-sectional estimates for consecutive time points will increase as the level (stage of sampling) and proportion of overlap in the sampled units increases. Also, the degree of increase in this correlation will vary from characteristic to characteristic.

Prior to 1988, there was no sample overlap at the household level for SCF between two consecutive years. However, there always has been a considerable, but not complete, PSU overlap between the SCF samples for two consecutive years. In order to evaluate the impact of household overlap on the estimation and/or the detection of change in average income, it was decided to introduce a one rotation group sample overlap at the household level for the 1988 and 1989 survey years.

The rotation groups 1, 2, 3, 5 and 6 of the 1988 LFS April sample were used for the 1988 SCF. The rotation groups 1, 2, 3, and 6 of the 1989 LFS April

sample and rotation group 5 of the 1988 LFS April sample were used for the 1989 SCF. Thus rotation group 5 of the 1988 LFS sample was used in both the 1988 and the 1989 surveys. In practice, however, only the 1988 SCF respondents of this group were approached for data collection in 1989, i.e. nonrespondents in 1988 were not contacted again. In the case of overlapping PSU's, there will be considerable enumeration area overlap with the selected households one year apart being neighbours.

In this study, by using the data from the 1988 and 1989 SCF's, we evaluate the relative impact of PSU and household sample overlap on the estimation and/or the detection of change in average income (in current dollars) for economic families of size 2-plus (E.F. 2+) and of unattached individuals, at provincial, regional, and national levels. Specifically, we attempt to answer the following two questions:

(i) Does household overlap provide estimates of change that are more efficient than those based on PSU overlap?

(ii) Should the t-statistic for testing the statistical significance of change in average income be based on the assumption of zero correlation or on the use of estimated correlation between two cross-sectional estimates?

In section 2, we introduce the necessary notation. The assumptions required to undertake this evaluation are also described. The need to make these assumptions, instead of strictly adhering to a particular design, system or methodology, is also dictated by the lack of data available for certain cases. The relative efficiencies for the estimation of change for the two types of overlap are examined in section 3. Section 4 describes the impact of using the

estimated value of correlation for statistical inference about change of income based on a t-test. Section 5 briefly discusses other issues, such as the effect of household overlap on cross-sectional estimates and on longitudinal data. The main recommendations of the report are also included in that section.

2. NOTATION

We will use two sample design structures for this study. These are:

D_1 - The sample for this design is composed of rotation groups 1,2,3 and 6 of the 1988 and 1989 SCF's. Note that the 1988 and 1989 components have PSU type overlap but no household overlap as rotation group 5 is excluded.

D_2 - The sample for this design is composed of rotation group 5 of the 1988 LFS (used in the 1988 and 1989 SCF's) and three out of four rotation groups (1,2,3, or 6) for the 1988 and 1989 SCF's. Note that: (i) under this design the 1988 and 1989 SCF components theoretically have a 25% overlap at the household level, and (ii) we are not using the full samples of the 1988 and 1989 SCF's.

The reason for using a smaller sample for design D_2 is so that various comparisons between D_1 and D_2 will be based on samples of approximately equal size. Also four different samples can be generated from the available data that conform to design D_2 and hence can generate relatively more stable estimates of some quantities (see Appendix I for details). It should be pointed out that the income, reported for SCF 1988, is that for the previous year, i.e., 1987. A similar observation also applies to the SCF 1989. For $i=1,2$, let

$\hat{X}_{88(i)}$ - estimate of 1988 average income X_{88} (in constant dollars) of an economic unit based on sample D_i , and

$\hat{V}(\hat{X}_{88(i)})$ - estimate of variance for $\hat{X}_{88(i)}$.

The corresponding estimates for the income year 1987 will use the subscript "87" in place of "88". In addition, we define

$\hat{\rho}_4$ - estimate of correlation coefficient ($\hat{X}_{88(1)}, \hat{X}_{87(1)}$),

$\hat{\gamma}_4$ - estimate of correlation coefficient ($\hat{X}_{88(2)}, \hat{X}_{87(2)}$)

Note that $\hat{\rho}_4$ is based on a D_1 type sample whereas $\hat{\gamma}_4$ is based on D_2 type sample. Some mathematical details about computing variances and correlation coefficients are given in Appendix I.

Let

\hat{d}_1 - estimate of change in average income, $X_{88} - X_{87}$, based on sample D_1 ,

- $\hat{X}_{88(1)} - \hat{X}_{87(1)}$, and

\hat{V}_1 - estimate of variance for \hat{d}_1 .

A mathematical summary of the Jackknife method for estimating the variance of change is given in Appendix II.

The relative efficiency, R , of design D_2 relative to design D_1 is given by the equation:

$$R = \frac{1/V_2}{1/V_1} * 100$$

$$= \frac{V_1}{V_2} * 100 \quad (2.1)$$

where V_1 and V_2 are the variances of estimate of change \hat{d}_1 and \hat{d}_2 , under designs D_1 and D_2 , respectively. Since these V 's are unknown, they will be replaced by the corresponding estimates (\hat{V} 's) in computing \hat{R} (the estimate of R). We make the following observations regarding the interpretation of \hat{R} :

- (i) if $\hat{R} > 100$ then the design D_2 is said to be more efficient than the design D_1 ,
- (ii) if $\hat{R} < 100$ then the design D_2 is said to be less efficient than the design D_1 , and
- (iii) if $\hat{R} = 100$ then the two designs are equally efficient. However, if \hat{R} is close to 100 then the two designs are almost equally efficient.

In our context, $\hat{R} > 100$ means that the household overlap provides better estimates of change than those provided by PSU overlap. How much better depends on how far \hat{R} is from the value of 100. Similar observations can be made about other values of \hat{R} .

3. RELATIVE EFFICIENCIES OF TWO DESIGNS

The relative efficiencies for estimates of change in average income for E.F. 2+ and unattached individuals have been computed and are presented in Tables I and II respectively. It should be noted that these relative efficiencies are for average income estimates in constant dollars.

3.1 Average Income of E.F. 2+

The $\hat{\rho}_4$, $\hat{\gamma}_4$ and relative efficiency estimates for E.F. 2+ are presented in Table I at the province, regional and Canada levels. As expected, the estimate $\hat{\gamma}_4$ is generally larger than $\hat{\rho}_4$. This is true for the Prairie region and Canada level estimates. At the province level, the $\hat{\gamma}_4$ is greater than $\hat{\rho}_4$ for all provinces except New Brunswick and Alberta. For those two provinces, $\hat{\gamma}_4$ is slightly smaller than $\hat{\rho}_4$. These two situations can be attributed to the sampling variability as we are looking at the estimates of the correlations and not at the

actual correlations.

The relative efficiencies, calculated by using $\hat{\rho}_4$ and $\hat{\gamma}_4$, show that the gains due to 25% household overlap are marginal at the Prairie region and Canada levels. The gains are in the range from 4% to 6% for those two areas. The gains due to household overlap at the province level are between 1% and 10% for six provinces, in the 15-20% range for Manitoba and Saskatchewan, and 40% for P.E.I.

Due to the fact that $\hat{\rho}_4$ is slightly larger but very close to $\hat{\gamma}_4$ for New Brunswick and Alberta, the relative efficiencies for these two provinces are less than, but close to, 100%. For these two provinces, the 25% household overlap does not provide an improvement for the estimate of change.

These results indicate that there are some, but not large, gains in relative efficiency for estimates of change in average income for E.F. 2+ at the various levels due to household overlap. However, these gains are marginal for household overlap in relation to PSU overlap, and thus, 25% household overlap does not result in a significant improvement for estimates of change.

The relative efficiency, R can also be approximately estimated by $\frac{1 - \hat{\rho}_4}{1 - \hat{\gamma}_4}$ under the assumption of equal variances for the two years and two designs. The calculations (in Table I and Table II) agree with this assumption in most cases.

3.2 Average Income of Unattached Individuals

In Table II, relative efficiencies are presented for estimates of change in the average income of unattached individuals. Since 'unattached individuals' is a considerably smaller characteristic than 'E. F. 2+', one can expect more unstable variance estimates and correlation coefficient estimates for that characteristic.

The Prairie region and Canada level estimates show a gain, due to household overlap, in relative efficiency of 8% and 10% respectively; there is a small decline in relative efficiency for the Atlantic region. However, this decline can be attributed to the sampling variation of various estimates used in computing relative efficiencies. Most of the gains at the provincial level are between 6% and 24%. This range is consistent with the observations from Table I (E.F. 2+).

Generally speaking, for estimating the change in average income of unattached individuals, 25% household sample overlap is more efficient than PSU type overlap with the exceptions of Nova Scotia and Quebec. However, the efficiency gains are not substantial.

In summary, the relative efficiencies for estimating change in average income of both E.F. 2+ and unattached individuals indicates that the design with the household overlap is more efficient than the one with PSU overlap. However, the magnitudes of relative efficiencies do not clearly demonstrate substantial superiority of household overlap over PSU overlap with respect to the characteristics under study.

4. DETECTION OF CHANGE IN AVERAGE INCOMES

The question of testing the statistical significance of change in average income between two consecutive years at province level is of interest to various analysts. The conclusions about this question have usually been arrived at by using the appropriate t-statistic under the assumption that the corresponding estimates for two consecutive time points are independent or uncorrelated. We examine the impact on t-statistic values of use of the estimated correlation

coefficient based on the PSU overlap, i.e. $\hat{\rho}_4$, relative to the assumption of zero correlation. Specifically, we are interested in testing the hypothesis of no change in average income for an economic unit from 1987 to 1988 against the alternative of change in average income. The value 1.96, corresponding to a 5% level of significance, is usually used to infer about the change in average income.

Based on the empirical evidence about the correlation coefficient $\hat{\rho}_4$ from Tables I and II, it can be said that the estimates of average incomes for 1987 and 1988 are positively correlated. In other words, the variances for estimates of change under the PSU overlap (Design 1) are less than that under the assumption of zero correlation.

Let,

$$t_1 = \hat{d}_1 / (\hat{V}_1)^{1/2} \quad (4.1)$$

and let t_1^* be the t-value obtained from t_1 by taking $\hat{\rho}_4=0$.

Due to the positive correlation, i.e., $\hat{\rho}_4 > 0$, that is observed between the estimates of average income for the two years, it is easy to see that $t_1 > t_1^*$.

Analysis using a t-statistic to detect change in average income in constant dollars was carried out for two cases with different correlations. The two cases are: (i) $\hat{\rho}_4=0$; and (ii) $\hat{\rho}_4$ from Tables I or II. There are three possibilities regarding the changes in conclusion as to the significance or nonsignificance of the t-statistic. These are:

(a) The t-test is significant in both cases

Note that if $t_1^* > 1.96$, then $t_1 > 1.96$. Therefore a t-statistic

significant under the assumption of zero correlation is also significant under the use of estimated $\hat{\rho}_4$. That is, there is no difference in the conclusions under the two cases. This is true in the case of average income for E.F. 2+ at the Prince Edward Island and Canada levels (see Table III). For unattached individuals, this is also the case with New Brunswick (see Table IV).

(b) The t-test is not significant in both cases

If t_1 and t_1^* are both less than 1.96, then the t-statistic is significant in neither case. Again, there is no difference in conclusions under these two cases.

These include the situations in Newfoundland, Nova Scotia, New Brunswick, Quebec, Ontario, Manitoba, Saskatchewan, Alberta and British Columbia, as well as the two regional estimates for E.F. 2+. As far the tests for unattached individuals, conclusions for all provinces and regions are of insignificant change under both cases except in New Brunswick. In these cases, either the changes in average incomes are small or the correlation coefficients are small. This explains the statistical nonsignificance of changes in income, since it is difficult to detect significant change when changes are small.

(c) The t-test is not significant assuming zero correlation but significant using estimated $\hat{\rho}_4$

If $t_1^* < 1.96$ but $t_1 > 1.96$, then the conclusions are different in the two cases, i.e., the t-statistic for change of income is statistically insignificant under zero correlation but significant under estimated $\hat{\rho}_4$.

This case is of most interest as it results in different conclusions under the two situations. However, in our study, no such case occurred.

Another useful statistic, the P-value, is reported in Table III and Table IV. The P-value, or significance level, allows readers to set their own level of significance and to evaluate the results according to their own requirements.

In summary, for this study, the t-tests carried out using estimated $\hat{\rho}_4$ do not produce any changes in conclusions regarding the detection of change. In other words, since the estimated correlation coefficients are small, the conclusions based on use of $\hat{\rho}_4$ do not differ much from those based on the assumption of zero correlation.

5. CONCLUSIONS

As was pointed out in Section 3, the samples with a theoretical 25% household overlap do not substantially improve the estimates of change in average income. This is due to the high positive intra-cluster correlation for a characteristic such as average income. Sampling the same household compared to sampling two different households in the same cluster or PSU on two consecutive occasions does not significantly improve the estimates of change. In the 1988 and 1989 SCF's there was actually at most a 20% sample overlap at the household level. Thus, the corresponding relative efficiencies will be lower than those presented in tables I and II. As well, the information about clusters (1988 rotation group 5) was not updated to reflect any growth in those areas. The inability to update such information can adversely affect the cross-sectional estimates and hence, the estimates of change. Furthermore, we do not have the 1989 LFS data for rotation group 5 of 1988, such as the updated household membership composition. This restricts proper manipulation of the 1989 SCF data

based on 5 rotation groups with respect to the weighting and imputation.

The rotation pattern 4-8-4 (four months in the sample, eight months out of the sample and a final four months in the sample) used in the U.S. Current Population Survey has the advantage of exploiting a sample with 50% household overlap for the estimation of change in income without any adverse effects on the cross-sectional estimates as the updated lists are used and the corresponding labour force data are collected.

The usefulness of longitudinal data for various analytic purposes is well known. For example, it is useful for policy purposes to monitor the movement of economic units between "above" and "below" poverty lines. If household overlap is to be used, its principal focus should be guided not by improvement in estimates of change but by other uses of longitudinal data. If the focus is on improved estimation of change, then various operational steps or different LFS rotation patterns should be used to overcome the possible adverse effects of household overlap on cross-sectional estimates. With the present rotation design, these steps would require updating information about the clusters and getting current LFS data. Alternately, another rotation scheme for LFS that alleviates these problems, can be used.

Now we discuss the second question of using an estimate of the correlation $\hat{\rho}$ in computing t-statistics for inference about change. The use of complete SCF samples (including those non-matching PSUs) for estimating the correlation would have required considerable methodological and systems development. To overcome the need for this major development it was decided to use the current jackknife variance estimation system for this purpose. This system required the use of

PSU's that are common to the 1988 and 1989 SCF samples. Furthermore, we have only one estimate of this correlation coefficient, $\hat{\rho}_4$ (based on the 1988 and 1989 surveys), and for this reason are unable to conclude that $\hat{\rho}_4$ is always significantly different from zero; nor can we answer any questions about the stability of the estimate $\hat{\rho}_4$. It is suggested that we continue with the present practice of using $\hat{\rho}=0$ in computing t-statistics until methodological research can produce theoretically sound correlation estimates.

TABLE I

CORRELATION COEFFICIENTS AND RELATIVE EFFICIENCIES FOR E.F. 2+

PROVINCE	$\hat{\rho}_4$	$\hat{\gamma}_4$	\hat{R}
NEWFOUNDLAND	0.16	0.20	104
PRINCE EDWARD ISLAND	0.26	0.49	140
NOVA SCOTIA	0.34	0.35	101
NEW BRUNSWICK	0.24	0.22	97
ATLANTIC REGION	0.29	0.29	100
QUEBEC	0.35	0.37	103
ONTARIO	0.10	0.14	105
MANITOBA	0.16	0.28	115
SASKATCHEWAN	0.20	0.33	120
ALBERTA	0.37	0.35	97
PRAIRIE REGION	0.29	0.31	104
BRITISH COLUMBIA	0.13	0.15	103
CANADA	0.16	0.20	106

TABLE II

CORRELATION COEFFICIENTS AND RELATIVE EFFICIENCIES FOR UNATTACHED INDIVIDUALS

PROVINCE	$\hat{\rho}_4$	$\hat{\gamma}_4$	\hat{R}
NEWFOUNDLAND	0.00	0.04	106
PRINCE EDWARD ISLAND	0.00	0.11	116
NOVA SCOTIA	0.47	0.26	73
NEW BRUNSWICK	0.13	0.14	102
ATLANTIC REGION	0.14	0.11	96
QUEBEC	0.31	0.23	90
ONTARIO	0.15	0.26	116
MANITOBA	0.00	0.16	117
SASKATCHEWAN	0.20	0.27	110
ALBERTA	0.15	0.20	107
PRAIRIE REGION	0.12	0.19	108
BRITISH COLUMBIA	0.00	0.19	124
CANADA	0.14	0.22	110

TABLE IIITABLE OF 't-VALUES' AND 'P-VALUES' FOR E.F. 2+

PROVINCE	t_1	t_1^*	t-test result	P_1	P_1^*
NFLD.	1.81	1.66	nc	0.070	0.097
P.E.I.	3.26	2.80	nc	0.001	0.005
N.S.	0.31	0.25	nc	0.757	0.803
N.B.	0.91	0.79	nc	0.363	0.430
ATLANTIC	0.70	0.59	nc	0.484	0.555
QUE.	0.48	0.39	nc	0.631	0.697
ONT.	1.92	1.83	nc	0.055	0.194
MAN.	1.74	1.61	nc	0.082	0.107
SASK.	0.45	0.40	nc	0.653	0.689
ALBT.	0.65	0.51	nc	0.516	0.610
PRAIRIES	1.30	1.09	nc	0.194	0.276
B.C.	1.00	0.93	nc	0.317	0.352
CANADA	2.23	2.05	nc	0.026	0.040

- Notes :
- 1) t_1 and t_1^* are the t-values based on estimated $\hat{\rho}_4$ and on $\hat{\rho}_4 = 0$.
 - 2) 'nc' and 'c' represent 'no change' and 'change' in the conclusion under the two assumptions about $\hat{\rho}_4$.
 - 3) P_1 and P_1^* are the P-values of t_1 and t_1^* , respectively.
 - 4) t-values and P-values are based on average incomes in constant dollars.

TABLE IV

TABLE OF 't-VALUES' AND 'P-VALUES' FOR UNATTACHED INDIVIDUALS

PROVINCE	t_1	t_1^*	t-test result	P_1	P_1^*
NFLD.	0.87	0.87	nc	0.384	0.384
P. E. I.	0.24	0.24	nc	0.810	0.810
N. S.	1.18	0.86	nc	0.238	0.390
N. B.	3.39	3.18	nc	0.001	0.001
ATLANTIC	1.07	1.00	nc	0.285	0.317
QUE.	1.42	1.18	nc	0.156	0.238
ONT.	0.52	0.48	nc	0.603	0.631
MAN.	0.00	0.00	nc	1.000	1.000
SASK.	0.30	0.27	nc	0.764	0.787
ALBT.	0.13	0.12	nc	0.897	0.904
PRAIRIES	0.30	0.28	nc	0.764	0.779
B. C.	1.19	1.19	nc	0.234	0.234
CANADA	0.12	0.11	nc	0.904	0.912

* The notes for Table III are applicable to this table as well.

APPENDIX I

Estimation of Correlation Coefficients

We provide some details about the estimation of the correlation coefficients $\hat{\rho}_4$ and $\hat{\gamma}_4$. It was mentioned in Section 2 that two designs, D_1 and D_2 , would be compared. The methodology for variance estimation will be described in Appendix II. Five samples or ten sub-samples were used to generate various estimates. These are:

Sample No.	1988 Survey Sub-sample A Rotation Groups	1989 Survey Sub-sample B Rotation Groups
1	1,2,3,6	1,2,3,6
2	1,2,3,5	1,2,3,5
3	1,2,5,6	1,2,5,6
4	1,3,5,6	1,3,5,6
5	2,3,5,6	2,3,5,6

For each sample, only the PSU's that are common to sub-samples A and B are retained. Note that the A and B components of the last four samples have rotation group 5 in common, thereby allowing the two sub-samples to have 25% household overlap, whereas sample 1 has PSU overlap only.

For $i=1,2,\dots,5$, and $y=87,88$, let:

$\hat{V}_i(\hat{X}_y)$ - variance estimate of \hat{X}_y based on the appropriate sub-sample of the i^{th} sample, and

$\hat{V}_i(\hat{X}_{88}-\hat{X}_{87})$ - variance estimate of $\hat{X}_{88} - \hat{X}_{87}$ based on the i^{th} sample.

Thus we have:

$$\hat{\rho}_4 = \frac{\hat{V}_1(\hat{X}_{88}) + \hat{V}_1(\hat{X}_{87}) - \hat{V}_1(\hat{X}_{88} - \hat{X}_{87})}{2[\hat{V}_1(\hat{X}_{88})\hat{V}_1(\hat{X}_{87})]^{1/2}} \quad (\text{A1.1})$$

and

$$\hat{\gamma}_4 = \frac{1}{4} \sum_{i=2}^5 \frac{\hat{V}_1(\hat{X}_{88}) + \hat{V}_1(\hat{X}_{87}) - \hat{V}_1(\hat{X}_{88} - \hat{X}_{87})}{2[\hat{V}_1(\hat{X}_{88})\hat{V}_1(\hat{X}_{87})]^{1/2}} \quad (\text{A1.2})$$

Note that $\hat{\rho}_4$ is based on one sample only, i.e., sample 1, while $\hat{\gamma}_4$ is based on four samples, i.e., samples 2, 3, 4 and 5. When computing relative efficiencies of change, the estimate of variance of

$$\frac{1}{4} \sum_{i=1}^4 \hat{V}_1(\hat{X}_y) \text{ , where } y=87,88, \text{ is used for } \hat{V}(\hat{X}_y).$$

APPENDIX II

Estimation of the Change and Variance

The estimator of characteristic X is the regression estimator, which used an integrated weighting method, see [2] for more details. It is of the form

$$\hat{X} = (X_t' \Pi^{-1} Z) (Z' \Pi^{-1} Z)^{-1} P \quad (A2.1)$$

where X_t is a transformed matrix of the characteristic of interest for each unit of the target population, Π is a matrix of weights corresponding to each unit of the target population, Z is a transformed matrix of auxiliary variables (age/sex groups, for example) defined for each unit of the target population and P is a matrix of the population control totals for each unit of the target population. The dimensions of above matrices depend on the number of units defined for the target population.

The variance estimate of \hat{X} is the Jackknife variance estimate given by

$$\text{Var}(\hat{X}) = \sum_h \frac{(n_h - 1)}{n_h} \sum_i (\hat{X}_{hi} - \hat{X})^2 \quad (A2.2)$$

where n_h is the number of PSUs (replicates) in stratum h , \hat{X}_{hi} is the Jackknife estimate for each replicate i of stratum h , calculated by (A2.1), and \hat{X} is the full sample (no removal of replicates) estimate calculated by (A2.1).

As previously mentioned, \hat{d} is the estimate of change for the value \hat{X} from income year 1987 to 1988. The estimate of this change is given by

$$\hat{d} = \hat{X}_{88} - \hat{X}_{87}.$$

The variance of \hat{d} is estimated by

$$\begin{aligned} \text{Var}(\hat{d}) &= \text{Var}(\hat{X}_{88} - \hat{X}_{87}) \\ &= \sum_h \frac{(n_h - 1)}{n_h} \sum_i [(\hat{X}_{hi88} - \hat{X}_{hi87}) - (\hat{X}_{88} - \hat{X}_{87})]^2 \end{aligned} \quad (A2.3)$$

where we included only those PSU's that are common to survey years 1988 and 1989.

The variance of change with the assumption of zero correlation used in Section 4 is

$$\begin{aligned}\text{Var}(\hat{d}) &= \text{Var}(\hat{X}_{88} - \hat{X}_{87}) \\ &= \text{Var}(\hat{X}_{88}) + \text{Var}(\hat{X}_{87})\end{aligned}\tag{A2.4}$$

where the variance formula (A2.2) is applied individually to each year.



1010072836

Ca OOS

ACKNOWLEDGEMENTS

The authors would like to thank G. Gray for his insightful discussions during the progress of the work and also for his valuable comments on an earlier draft of the manuscript.

REFERENCES

1. Kish, L. (1965), Survey Sampling, John Wiley & Sons, New York.
2. Lemaitre, G. and Dufour, J. (1987), "An Integrated Method for Weighting Persons and Families". Survey Methodology, Vol.13, No.2. Statistics Canada.
3. Statistics Canada (1990), Methodology of the Canadian Labour Force Survey.

