

Catalogue no. 62M0004XCB

# **User Guide for the Public-use Microdata File**

## **Survey of Household Spending, 2006**

July 2008

**Income Statistics Division  
Statistics Canada, Ottawa, K1A 0T6  
Telephone: 613 951-7355**

*Ce document est disponible en français.*

“Income Statistics Division, Statistics Canada” must be credited when reproducing or quoting any part of this document.

# Table of contents

<b>Introduction .....</b>	<b>3</b>
Background.....	3
New for 2006.....	3
Other documents.....	4
For further information.....	4
<b>Technical characteristics of the file.....</b>	<b>6</b>
<b>Survey methodology.....</b>	<b>7</b>
The survey universe.....	7
Survey content and reference period.....	7
The sample .....	8
Data collection.....	8
Data processing and quality control.....	8
Weighting, re-weighting, and Census historical revision of SHS.....	9
<b>Data quality.....</b>	<b>10</b>
Sampling error.....	10
Non-sampling error.....	11
The effect of large values .....	14
Comparability over time.....	14
<b>Guidelines for tabulation, analysis and dissemination .....</b>	<b>15</b>
Guidelines for rounding .....	15
Guidelines for the weighting of the sample for totalling purposes .....	16
Types of estimates: categorical versus quantitative.....	16
Confidentiality of the public-use microdata .....	30
<b>Appendices—See accompanying Excel file.....</b>	<b>31</b>
Appendix A Frequency counts.....	31
Appendix B Averages, aggregates, minimum and maximum values.....	31
Appendix C Inclusion of spending variables in past microdata files .....	31
Appendix D Coefficients of variation for published data from the 2006 SHS.....	31

## Introduction

### Background

This public-use microdata file presents data from the 2006 Survey of Household Spending (SHS) conducted in January through April 2007. Information about the spending habits, dwelling characteristics and household equipment of Canadian households during 2006 was obtained by asking people in the 10 provinces to recall their expenditures for the previous calendar year (spending habits) or as of the time of the interview (dwelling characteristics and household equipment).

Conducted since 1997, the Survey of Household Spending integrates most of the content found in the Family Expenditure Survey and the Household Facilities and Equipment Survey. Many data from these two surveys are comparable to the Survey of Household Spending data. However, some differences related to methodology, to data quality and to definitions must be considered before comparing these data. See For further information below.

### New for 2006

Beginning with the 2006 SHS, the reference date for household composition, tenure, dwelling characteristics and household equipment are as of the time of the interview instead of December 31<sup>st</sup> of the reference year. The distinction between full-year and part-year members has been removed. Spending data are collected for all members of the household at the time of the interview.

For the 2006 reference year, automatic edits built into the electronic questionnaire replaced the balance edit and regional office editing performed in previous years.

The variable DVDPLYR in 2005 was renamed DVD in 2006.

The variable J019TOT in 2005 was renamed to J041 in 2006.

The following variables were deleted in 2006.

Operating a farm	OPFARM
The number of floors in a dwelling	NUMFLR
The type of dwelling the reference person lived previously	RPPRDWTY
Number of floors in the previous dwelling	RPPREFLR
Moved to a larger dwelling	LARGEDWG
Moved to a smaller dwelling	SMALLDWG
Moved to a less expensive dwelling	CHEAPDWG
Moved to a better quality dwelling or neighbourhood	BETTRDWG
Moved to be closer to facilities and services	CLOSEFAC
Moved to establish own household	ESTHLD
Moved to change from owner to renter or renter to owner	CHNGTEN
Moved for a new job or job transfer	CHNGJOB
Moved to be closer to work or school	CLOSWORK

Moved for family reasons	FAMREA
Moved for health reasons	HEALTHR
Moved for other reasons	OTHERR
Rent calculated based on income	RETOINC
Total number of person weeks member of the household	NUMWKSP

Number of persons members at some time in the year (HHSZTOTP) redefined to be number of persons at the time of the interview.

## Other documents

- Data dictionary (variable specifications, code sets and other information) is available in pdf format.
- Record layout is available in Excel format.
- Appendices are available in Excel format.
  - Appendix A presents the frequency counts for non-dollar variables in the public-use microdata file. They are included to help you verify your tabulations.
  - Appendix B presents expenditure data tabulated using the public-use microdata file and also using the internal survey database. They are included to help you verify your tabulations.
  - Appendix C contains a table indicating the spending variables included in previous public-use microdata files of the Survey of Household Spending and the Family Expenditure Survey.
  - Appendix D shows any changes in variables from the previous years.
  - Appendix E presents the coefficients of variation for published data from the 2006 SHS.

## For further information

- Additional information about the SHS can now be obtained free on the Statistics Canada web site ([www.statcan.ca](http://www.statcan.ca)). See especially:
- Note to former users of data from the Family Expenditure Survey (62F0026MIE2000002)
- Note to former users of data from the Household Facilities and Equipment Survey (62F0026MIE2000003)
- User Guide for the Survey of Household Spending, 2006 (62F0026MIE2008001)

- Methodology for the Survey of Household Spending (62F0026MIE2001003)
- 2003 Survey of Household Spending Data Quality Indicators (62F0026MIE2005006)

For more information about the current survey results and related products and services, or to enquire about the concepts, methods or data quality of the Survey of Household Spending, contact Client Services (613-951-7355; 1-888-297-7355; fax 613-951-3012; *income@statcan.ca*), Income Statistics Division.

## Technical characteristics of the file

**Content:** Household spending, dwelling characteristics, and household equipment, 2006

**Source:** Survey of Household Spending, 2006  
Income Statistics Division  
Statistics Canada

### Data set definition:

**Data set name** ..... SHS2006.TXT

**Number of records** ..... 14,635

### Format

Record length ..... 2,066

## Survey methodology

(For more detailed information, see the *Methodology of the Survey of Household Spending* available free on the Statistics Canada web site at [www.statcan.ca](http://www.statcan.ca)).

### The survey universe

The 2006 Survey of Household Spending was carried out in private households in Canada's 10 provinces.<sup>1</sup> The territories were not included in 2006, being surveyed every odd numbered year.

The following groups were excluded from the survey:

- those living on Indian reserves and crown lands (with the exception of the territories);
- official representatives of foreign countries living in Canada and their families;
- members of religious and other communal colonies;
- members of the Canadian Forces living in military camps; and
- people living full time in institutions: for example, inmates of penal institutions and chronic care patients living in hospitals and nursing homes.

The survey covers about 98% of the population in the 10 provinces.

Spending data were collected for every household member at the time of the interview, including those who joined the household in 2006 or 2007 regardless of whether the previous household existed or the person was living alone. Data were not collected for those who left the household in 2006 or 2007. As a result, an important difference between the 2006 SHS and previous SHS methodology is the elimination of the distinction between "part-year" and "full-year" members and households.

Persons temporarily living away from their families (for example, students at university) were included in the household to avoid double counting.

### Survey content and reference period

Detailed information was collected about expenditures for consumer goods and services, changes in assets, mortgages and other loans, and annual income. This information was collected for the calendar year 2006 (the survey reference year). Information was also collected about dwelling characteristics (e.g., type and age of heating equipment) and household equipment (e.g., appliances, communications equipment, and vehicles). This type of information was collected as of the time of the interview.

Because the Survey of Household Spending is designed principally to provide detailed information on non-food expenditures, only an overall estimate of food expenditure is recorded. Detailed information on food expenditure is provided by

---

1. In order to reduce response burden for northern households, the Survey of Household Spending is conducted in the north only every second year, starting with 1999.

the Food Expenditure Survey, which is conducted every four to six years. It was last conducted in 2001. In February 2003, the results were published in Food Expenditure in Canada, 2001, Catalogue no. 62-554-XIE.

## **The sample**

The sample size for the 2006 Survey of Household Spending was 20,436 eligible households.

The regular SHS sample was a stratified, multi-stage sample selected from the Labour Force Survey (LFS) sampling frame. Sample selection comprised two main steps: the selection of clusters (small geographic areas) from the LFS frame and the selection of dwellings within these selected clusters. The LFS sampling frame mainly uses 2001 Census geography and 2001 population counts.<sup>2</sup>

## **Data collection**

The 2006 Survey of Household Spending was conducted from January to April 2007. Data were collected by computer assisted personal interview (CAPI) using a laptop personal computer. A copy of this questionnaire is available on request.

## **Data processing and quality control**

Due to introduction of the new electronic CAPI questionnaire, changes were made to the data processing and quality steps. Automatic edits built into the electronic questionnaire replace the balanced edit and regional office edits performed in previous years.

For the 2006 Survey of Household Spending, the interviewers recorded the information provided by the respondents using a laptop and performed the initial editing at the same time. For example, the range edit provided a minimum and maximum amount for certain purchases and was triggered if the amount entered by the interviewer was unusual. Other edits indicated inconsistencies in responses, e.g. if the household tenure was “renter” but no rent was paid.

The next stage of editing was done in the head office to verify unusual or high values and inconsistencies, and to correct invalid responses.

If a household indicated that it had an expense but could not provide the amount, these missing responses were imputed using the nearest neighbour method. Statistics Canada’s Canadian Census Edit and Imputation System (CANCEIS) were used to insert values from donor records having similar characteristics, chosen specifically to fit the variable. For example, total household income was used for most variables; dwelling type, household size and province were also frequently used.

---

2. A detailed description of the Labour Force Survey sampling frame can be found in *Methodology of the Canadian Labour Force Survey*, Statistics Canada, Catalogue no. 71-526-XIE.

Tabulation for the 2006 Survey of Household Spending was completed using a PC/client server-based system. This system provides tools (database querying, searching, and viewing capabilities) for spotting systematic errors.

## **Weighting, re-weighting, and Census historical revision of SHS**

Users should note that the weights for the SHS reference years 1997 to 2003 have been revised. These revisions were published along with the 2005 survey results in December 2006.

The estimation of population characteristics from a sample survey is based on the idea that each sampled household represents a certain number of other households in addition to itself. These numbers are called the survey weights of the sample. To improve the representativity of the sample, the weights are adjusted so that the estimates from the sample are in line with population totals, or benchmarks, from other independent sources of information that are considered reliable. This is called weight calibration.

SHS uses two sources for calibration. The first source is the Census of Population which provides demographic benchmarks. From 1997 to 2003, SHS used benchmarks derived from the 1996 Census. Since the Census is conducted once every five years, Statistics Canada projects the Census results for later years (up to the present), and then revises those estimates when the next Census data become available. The projections use a variety of secondary information, including administrative data on births, deaths and migration.

The second source used for adjusting the survey weights for SHS are T4 data from Canada Revenue Agency, which ensures that the estimated distribution of earners in the survey matches the one in the Canadian population.

It was decided to take advantage of this historical revision to also introduce an improved calibration strategy for the SHS weights. Improvements to the calibration strategy were deemed necessary to put emphasis on SHS needs (such as the age groups used for calibration) and to take into account the quality of the benchmarks. It was also felt that there were too many benchmarks leading to too many constraints on the weights, and that this produced undesirable results, such as negative weights, which were not acceptable.

The current calibration strategy is as follows:

- Age
  - At the provincial level there are controls for 8 age groups (0-6, 7-17, 18-24, 25-34, 35-44, 45-54, 55-64, 65+).
  - At the CMA level: two age groups (0-17, 18+)
- There are controls for three size of household categories (one person, two persons, 3+)

- T4 adjustments are made to the weights of the population for income from wages and salaries (0-25<sup>th</sup> percentile, 25<sup>th</sup>-50<sup>th</sup>, 50<sup>th</sup>-65<sup>th</sup>, 65<sup>th</sup>-75<sup>th</sup>, 75<sup>th</sup>-95<sup>th</sup>, 95<sup>th</sup>-100<sup>th</sup>)

The weights and calibration strategy were implemented for SHS for the years 1997 and onward resulting in revised estimates of household spending for each year up to 2003. Users of SHS data should take care to make comparisons using the re-weighted data.

## Data quality

(For more detailed information, see the *Survey of Household Spending Data Quality Indicators*, soon to be available free on the Statistics Canada web site at [www.statcan.ca](http://www.statcan.ca).)

## Sampling error

Sampling errors occur because inferences about the entire population are based on information obtained from only a sample of the population. The sample design, the variability of the data, and the sample size determine the size of the sampling error. In addition, for a given sample design, different methods of estimation will result in different sampling errors.

The design for the 2006 Survey of Household Spending was a stratified multi-stage sampling scheme. The sampling errors for multi-stage sampling are usually higher than for a simple random sample of the same size. However, the operational advantages outweigh this disadvantage, and the fact that the sample is also stratified improves the precision of estimates.

Data variability is the difference between members of the population with respect to spending on a specific item or the presence of a specific dwelling characteristic or piece of household equipment. In general, the greater these differences are, the larger the sampling error will be. In addition, the larger the sample size, the smaller the sampling error.

## Standard error and coefficient of variation

A common measure of sampling error is the standard error (SE). Standard error is the degree of variation in the estimates as a result of selecting one particular sample rather than another of the same size and design. It has been shown that the 'true' value of the characteristic of interest lies within a range of +/- 1 standard error of the estimate for 68% of all samples, and +/- 2 standard errors for 95% of all samples.

The coefficient of variation (CV) is the standard error expressed as a percentage of the estimate. It is used to indicate the degree of uncertainty associated with an estimate. For example, if the estimate of the number of households having a given dwelling characteristic is 10,000 households, and the corresponding CV is

5%, then the true value is between 9,500 and 10,500 households, 68% of the time and between 9,000 and 11,000 households, 95% of the time.

Standard errors for the 2006 Survey of Household Spending were estimated using the 'bootstrap' method. This method is suitable for variance estimation of non-smooth statistics such as quintiles. For more information on standard errors and coefficients of variation, refer to the Statistics Canada publication, Methodology of the Canadian Labour Force Survey, Catalogue no. 71-526-XIE.

Coefficients of variation are available on request (contact Client Services, Income Statistics Division, 1-888-297-7355; [income@statcan.ca](mailto:income@statcan.ca)).

### **Data Suppression**

For reliability reasons, estimates with CVs greater than 33% are normally suppressed. Since CVs are not calculated for all estimates, data suppression for the Survey of Household Spending has been based on a relationship between the CV and the number of households reporting expenditure on an item. Analysis of past survey results indicates that CVs usually reach this level when the number of households reporting an item drops to about 30. Therefore, data have been suppressed for spending on items reported by fewer than 30 households.

However, data for suppressed items do contribute to summary level variables. For example, the expenditure for a particular category of clothing might be suppressed but this amount forms part of the total expenditure estimate for clothing.

### **Non-sampling error**

Non-sampling errors occur because certain factors make it difficult to obtain accurate responses or responses that retain their accuracy throughout processing. Unlike sampling error, non-sampling error is not readily quantified. Four sources of non-sampling error can be identified: coverage error, response error, non-response error, and processing error.

#### **Coverage error**

Coverage error results from inadequate representation of the intended population. This error may occur during sample design or selection, or during data collection and processing.

#### **Response error**

Response error may be due to many factors, including faulty design of the questionnaire, interviewers' or respondents' misinterpretation of questions, or respondents' faulty reporting.

Several features of the survey help respondents recall their expenditures as accurately as possible. First, the survey period is the calendar year because it is probably more clearly defined in people's minds than any other period of similar

length. Second, expenditure on food can be estimated as either weekly or monthly expenses depending on the respondent's purchasing habits. Third, expenses on smaller items purchased at regular intervals are usually estimated on the basis of amount and frequency of purchase. Purchases of large items (automobiles, for example) are recalled fairly easily, as are expenditures on rent, property taxes, and monthly payments on mortgages. However, even with these items, the accuracy of data depends on the respondent's ability to remember and willingness to consult records.

In the years prior to the 2006 SHS survey, a data quality control measure called the balance edit check was used. This measure identified the records where the expenditure reported was more than 20% different from the sum of income and net change of assets for a household. The interviewer or senior interviewer was instructed to attempt to collect additional information to try to balance the expenses with income and changes in assets within 15%. Unbalanced questionnaires (more than 20% ) at processing stage were deemed unusable and were not included in estimates. In the 2006 SHS, with the introduction of computer assisted interviews, the balance edit was not used at the collection stage. Instead a number of automatic edits flagged entries that seemed unusual or inconsistent. However, when the balancing was applied at the processing stage, the number of unbalanced questionnaires in the 2006 SHS increased significantly, from 546 questionnaires for the 2005 reference year to 4,300 or 29.4% of the 14,635 completed questionnaires for 2006.

Discarding as unusable this number of questionnaires would seriously risk biasing the results, so a careful analysis was done comparing the balanced and the unbalanced questionnaires. There were few significant differences in the average and percentage reporting of expenses between the balanced and the unbalanced questionnaires. Most of the difference lay in the income and the change of assets reported on the unbalanced responses. We concluded that we could include the unbalanced questionnaires in the estimates of expenses, but users should note that the quality of the income and change of assets may be lower than in previous years.

For the 2007 SHS the electronic questionnaire will be modified to re-introduce the balance edit feature and ensure identification and correction of out-of-balance records during the interview and collection, as in previous years.

### **Non-response error**

Non-response error occurs in sample surveys because not all potential respondents cooperate fully. The extent of non-response varies from partial non-response to total non-response.

Total non-response occurs when the interviewer is unable to contact the respondent, no member of the household is able to provide information, or the respondent refuses to participate in the survey. Total non-response is handled by adjusting the basic survey weight for responding households to compensate for non-responding households. For the 2006 Survey of Household Spending, the overall response rate was 71.6%. See Table 1 for provincial response rates.

In most cases, partial non-response occurs when the respondent does not understand or misinterprets a question, refuses to answer a question, or is unable to recall the requested information. Imputing missing values compensates for this partial non-response.

The importance of the non-response error is unknown but in general this error is significant when a group of people with particular characteristics in common refuse to cooperate and where those characteristics are important determinants of survey results.

**Table 1**  
**Response rates, Canada and provinces, 2006**

	Eligible households <sup>1</sup>	Non-contacts	Refusals	Usables	Balanced <sup>2</sup>	Response rate <sup>3</sup>
						%
Newfoundland and Labrador	1,753	105	293	1,331	925	75.9
Prince Edward Island	873	44	185	617	424	70.7
Nova Scotia	2,013	167	401	1,412	959	70.1
New Brunswick	1,774	124	295	1,313	898	74.0
Quebec	2,648	190	542	1,883	1,418	71.1
Ontario	3,097	209	654	2,156	1,521	69.6
Manitoba	1,957	91	358	1,472	1,236	75.2
Saskatchewan	1,898	84	357	1,430	987	75.3
Alberta	2,011	137	460	1,380	879	68.6
British Columbia	2,412	174	528	1,641	1,088	68.0
<b>Canada</b>	<b>20,436</b>	<b>1,325</b>	<b>4,073</b>	<b>14,635</b>	<b>10,335</b>	<b>71.6</b>

1. There is no longer a distinction between part-year and full-year households.

2. expenditures = income plus net change in assets +/- 20%

3. Usable/eligible\*100

### Processing error

Processing errors may occur in any of the data processing stages, for example, during data entry, editing, weighting, and tabulation. See Data processing and quality control for a description of the steps taken to reduce processing error.

## The effect of large values

For any sample, estimates can be affected by the presence or absence of extreme values from the population. These extreme values are most likely to arise from positively skewed populations. The nature of the subject matter of the SHS lends itself to such extreme values. Estimates of totals, averages and standard errors may be greatly influenced by the presence or absence of these extremes.

## Comparability over time

Conducted since 1997, the Survey of Household Spending integrates most of the content found in the Family Expenditure Survey and the Household Facilities and Equipment Survey. Many variables from these two surveys are comparable to those in the Survey of Household Spending. However, some differences related to the methodology, to data quality and to definitions must be considered before making comparisons.

For more information, refer to *Note to Former Users of Data from the Family Expenditure Survey*, Catalogue no. 62F0026MIE2000002 and *Note to Former Users of Data from the Household Facilities and Equipment Survey*, Catalogue no. 62F0026MIE2000003. Both documents are available free of charge on the Statistics Canada web site ([www.statcan.ca](http://www.statcan.ca)).

Historical data from the 1997 to the 2003 surveys of household spending have been re-weighted using the weighting methodology described in the section Weighting. Historical comparisons between data from those surveys and data from recent years of the Survey of Household Spending should generally be made with re-weighted data, although the differences between survey estimates from the old and new methodologies appear to be minimal at a summary level. Certain populations or variables, however, may be more strongly affected.

Starting with the 1997 Survey of Household Spending, 'Tenants' maintenance, repair and alterations' and 'Insurance premiums' were reduced by the proportion of rent charged to business. This may affect comparisons with data from previous years.

For the 2001 and 2005 reference years, extra questions were included for use in the weighting of the Consumer Price Index. This change may affect some historical comparisons. For example, in both 2005 and 2001, questions were added under 'Personal care' to collect extra information about hair care products, makeup, fragrances, deodorants and oral hygiene products. As a result of these extra questions, respondents may have given more precise information and the increase in the estimated expenditures for Personal care in 2001 and 2005 may have been caused by an improvement in respondent recall. The effect of additional questions on estimates is difficult to quantify. However, in 2002, when the extra questions were removed, the estimate for Personal care spending decreased again. Beginning in 2006, these questions will be left on every year.

The section of the questionnaire which covers “Repairs and improvements of owned principal residences” was extensively revised. From 1997 to 2003, this section had three broad questions: “Additions, renovations and other alterations”; “Replacement or new installation of built-in equipment, appliances and fixtures”; and “Repairs and maintenance”. Starting with the 2004 Survey of Household Spending, there were fourteen detailed questions and two columns, giving respondents the opportunity to split the costs for each question between “Repairs and maintenance” and “Improvements and alterations”.

## **Guidelines for tabulation, analysis and dissemination**

This section describes the guidelines that users should follow when totalling, analysing, publishing or releasing data taken from the public-use microdata file.

### **Guidelines for rounding**

To ensure that estimates from this microdata file intended for publication or any other type of release correspond to estimates that would be obtained by Statistics Canada, we strongly recommend that users comply with the following guidelines for rounding estimates.

- a) Estimates in the body of a statistical table must be rounded to the nearest hundredth using the traditional rounding technique, i.e., if the first or only number to be eliminated is between 0 and 4, the preceding number does not change. If the first or only number to be eliminated is between 5 and 9, the value of the last number to be retained increases by 1. For example, when using the traditional technique of rounding to the nearest hundredth, if the last two numbers are between 00 and 49, they are replaced by 00 and the preceding number (denoting hundredths) stays as is. If the last two numbers are between 50 and 99, they are replaced with 00 and the preceding number increased by 1.
- b) Total partial sub-totals and total sub-totals in statistical tables must be calculated using their unrounded corresponding components, then rounded in turn to the closest hundredth using the traditional rounding technique.
- c) Means, ratios, rates and percentages must be calculated using unrounded components (i.e., numerators and/or denominators), and then rounded to a decimal using the traditional rounding technique.
- d) Totals and differences in aggregates (or ratios) must be calculated using their corresponding unrounded components, then rounded to the nearest hundredth (or decimal place) using the traditional rounding technique.
- e) If, due to technical or other limitations, a technique other than traditional rounding is used, with the result that the estimates to be published or released differ in any form from the corresponding estimates that would be obtained by Statistics Canada using this microdata file, we strongly advise

users to indicate the reasons for the differences in the documents to be published or released.

- f) Unrounded estimates cannot under any circumstances be published or released in any way whatsoever by users. Unrounded estimates give the impression that they are much more precise than they actually are.

## **Guidelines for the weighting of the sample for totalling purposes**

The sample design used for the SHS is not self-weighted, meaning that the households in the sample do not all have the same sampling weight. To produce simple estimates, including standard statistical tables, users must use the appropriate sampling weight. Otherwise, the estimates calculated using the microdata files cannot be considered as representative of the observed population and will not correspond to those that would be obtained by Statistics Canada using this microdata file. See Weighting, re-weighting, and Census historical revision of SHS.

Users should also note that depending on the method they use to process the weight field, some software packages may not produce estimates that correspond exactly to those of Statistics Canada using this microdata file.

## **Types of estimates: categorical versus quantitative**

Before discussing how SHS data can be totalled and analysed, it is useful to describe the two main types of estimations that may be produced from the microdata file for the Survey of Household Spending.

### **Categorical estimates**

Categorical estimates are estimates of the number or percentage of households in the survey's target population that have certain characteristics or belong to a defined category. The number of households reporting a particular expenditure is an example of this type of estimate. The expression 'aggregate estimate' can also be used to refer to an estimate of the number of individuals with a given characteristic.

### **Examples of categorical questions:**

Does anyone in your household use the Internet from home?      \_yes    \_no

When was this dwelling originally built?

- 1945 or earlier
- 1946-1960
- 1961-1970
- 1971-1980
- 1981-1990
- 1991-2007

Is your dwelling:

- Owned without a mortgage by your household?
- Owned with (a) mortgage(s) by your household?
- Rented by your household?
- Occupied rent-free by your household?

### **Totalling of categorical estimates**

Estimates of the number of persons with a given characteristic can be obtained from the microdata file by adding the final weights of all records containing the desired characteristic or characteristics. Percentages and ratios in the X/Y form are obtained as follows:

- a) by adding the final weights of records containing the desired characteristic for the numerator X;
- b) by adding the final weights of records containing the desired characteristic for the denominator Y;
- c) by dividing the estimate for the numerator by the estimate for the denominator.

### **Quantitative estimates**

Quantitative estimates are estimates of totals or means, medians or other central tendency measurements of quantities based on all members of the observed population or based on some of them. They also explicitly include estimates in the form X/Y where X is an estimate of the total quantity for the observed population and Y is an estimate of the number of individuals in the observed population who contribute to that total quantity.

An example of a quantitative estimate is mean annual expenditure for personal and health care per household in the target population. The numerator corresponds to an estimate of total annual expenditure for personal and health care, and the denominator corresponds to an estimate of the number of households in the population.

### **Example of quantitative question:**

In 2006, how much did your household spend for telephone services? \_\_\_\_\_

### **Totalling of quantitative estimates**

Quantitative estimates can be obtained from the microdata file by multiplying the value of the desired variable by the final weight of each record, and then adding this quantity for all records of interest. For example, to obtain an estimate of total expenditure by households that were owners at the time of interview for electricity, the value reported for the question "In 2006, how much did your household spend on electricity?" is multiplied by the final weight of the record, and then that result is summed over all records with a positive response to the question "Is your house: 'Owned mortgage-free by your household' or 'Owned with one or more mortgages by your household'."

To obtain a weighted mean expressed by the formula  $X/Y$ , the numerator  $X$  is calculated as a quantitative estimate and the denominator  $Y$  as a categorical estimate. For example, to estimate mean household expenditures for electricity by owners, you must:

- a) estimate the total expenditure for electricity for households where the residence is owned, using the method described above;
- b) estimate the number of owned households by adding the final weights for all records with a positive response to the question "Is your house: 'Owned mortgage-free by your household' or 'Owned with one or more mortgages by your household'"; and then,
- c) divide the estimate obtained in a) by the one calculated in b).

### **Guidelines for statistical analysis**

The Survey of Household Spending is based on a complex survey design that includes stratification and multiple stages of selection, as well as uneven respondent selection probabilities. The use of data from such complex surveys poses problems for analysts, because the survey design and the selection probabilities influence the estimation and variance calculation methods to be used.

Although numerous analytical methods in statistical software packages allow for the use of weights, the meaning or definition of weights differs from that suitable for a sample survey. As a result, although the estimates done using those packages are in many cases accurate, **the variances calculated have almost no significance.**

For numerous analytical techniques (for example, linear regression, logistic regression, variance analysis), there is a way to make the application of standard packages more significant. If the weights of the records contained in the file are converted so that the mean weight is (1), the results produced by standard packages will be more reasonable and will take into account uneven selection probabilities, although they still cannot take into account the stratification and the cluster distribution of the sample. The conversion can be done using in the analysis a weight equal to the original weight divided by the mean of original weights for sampling units (households) that contribute to the estimator in question. However, because this method still does not take into account sample design stratification and clusters, the estimates of the variance calculated in this way will very likely be underestimates of true values.

### **Guidelines for release**

Before releasing and/or publishing estimates taken from the microdata file, users must first determine the level of reliability of the estimates. The quality of the data is affected by the sampling error and the non-sampling error as described above. However, the level of reliability of estimates is determined solely on the basis of sampling error, as evaluated using the coefficient of variation (CV) as shown in

the table below. In addition to calculating CVs, users should also read the section of this document regarding the characteristics of data quality.

Whatever CV is obtained for an estimate from this microdata file, users should determine the number of sampled respondents who contribute to the calculation of the estimate. If this number is less than 30, the weighted estimate should not be released regardless of the value of the CV for this estimate. For weighted estimates based on sample sizes of 30 or more, users should determine the CV of the rounded estimate following the guidelines below.

**Figure 2**  
**Sampling variability guidelines**

Type of Estimate	CV (in %)	Guidelines
1. Acceptable	0.0 – 16.5	Estimates can be considered for general unrestricted release. Requires no special notation.
2. Marginal	16.6 – 33.3	Estimates can be considered for general unrestricted release but should be accompanied by a warning cautioning subsequent users of the high sampling variability associated with the estimates. Such estimates should be identified by the letter M (or in some other similar fashion).
3. Unacceptable	Greater than 33.3	<p>Statistics Canada does not recommend the release of estimates of unacceptable quality. However, if the user chooses to do so then estimates should be flagged with the letter U (or in some other similar fashion) and the following warning should accompany the estimates:</p> <p>“The user is advised that . . . (specify the data) . . . do not meet Statistics Canada’s quality standards for this statistical program. Conclusions based on these data will be unreliable and most likely invalid.”</p>

### Computation of approximate CVs

In order to provide a way of assessing the quality of estimates, Statistics Canada has produced a coefficient of variation table (CV table) which is applicable to estimates of averages, ratios and totals obtained from this public use microdata file for the major variables of the SHS by province and at the Canada level (see Appendix E). The CV of an estimate is defined to be the square root of the variance of the estimate divided by the estimate itself and expressed as a percentage. The numerator of the CV is a measure of the sampling error of the

estimate, called the standard error, and is calculated at Statistics Canada with the bootstrap method. This method requires, among other things, information about the strata and the clusters, which can't be given on the public use microdata file for reasons of confidentiality. So that users may estimate CVs for variables not included in the CV tables, Statistics Canada has produced a set of rules to obtain approximate CVs for a wide variety of estimates. It should be noted that these rules provide approximate and, therefore, unofficial CVs. The quality of the approximation, however, is quite satisfactory, especially for the most reliable estimates. Note that accuracy of this approximation is reduced when the domains become smaller. Therefore, the CV approximation method must be used prudently when the domains are small. The document on data quality for the 1997 SHS contains the results of the evaluation of the performance of the CV approximation method.

### **How to obtain approximate CVs**

The following rules should enable the user to determine the approximate coefficients of variation for estimates of totals, means or proportions, ratios and differences between such estimates for sub-populations (domains) for which the Bootstrap CV is not provided in the CV tables.

**Important:** If the number of observations on which an estimate is based is less than 30, the weighted estimate should not be released regardless of the value of the CV for this estimate.

### **Rule 1: Approximating CVs for estimates of totals (aggregates)**

All the steps below must be followed to obtain an approximate CV (ACV) for an estimate of a total (either a number of households possessing a certain characteristic (categorical estimate) or a total of some expense for all households (quantitative estimate)) for a sub-population (domain) of interest:

- 1) Create a binary variable for each household, say  $I$ , equalling 1 if the household is part of the domain of interest, i.e. possesses the desired characteristic and 0 otherwise;
- 2) To estimate a quantitative variable, create a variable  $Y$  representing the product of the binary variable  $I$  and the variable of interest. To estimate a categorical variable, create a variable  $Z$  equal to 1 if the categorical variable is equal to the value of interest, and equal to 0 otherwise. Define variable  $Y$  as the product of  $I$  and  $Z$ ;
- 3) Do step (4) to step (9) for each province separately;
- 4) Calculate the sum over all the households of the product of the final weight (section Weighting), and  $Y$  (this sum represents the estimate of the total for the domain of interest in the province under consideration);
- 5) Calculate the sum over all the households of the product of the final weight and the household size;
- 6) Divide the result obtained in step (4) by the result obtained in step (5);
- 7) For each household, multiply the result obtained in step (6) by the household size;
- 8) For each household, define a variable, say  $E$ , by the subtraction of the result obtained in step (7) from  $Y$ ;

- 9) Calculate the sum over all the households of the product of the final weight minus 1, the final weight and E squared; (this sum represents the estimated variance of the total estimated at step 4);
- 10) Add up the result obtained in step (9) for each province;
- 11) The ACV is defined to be 100 times the square root of the result obtained in step (10), divided by the estimate. The estimate is the sum over all the provinces of the result obtained in step (4).

More formally, steps 1 to 10 above can be obtained with the following formula:

$$\sum_{p=1}^{11} \sum_{k \in S_p} (w_k - 1) w_k \left( Y_k - m_k \frac{\sum_{k \in S_p} w_k Y_k}{\sum_{k \in S_p} w_k m_k} \right)^2$$

where the index  $p$  corresponds to provinces,  $S_p$  is the sample of respondents for the province  $p$ , the index  $k$  corresponds to households,  $w_k$  is the final weight for the  $k^{\text{th}}$  household,  $m_k$  is the household size for the  $k^{\text{th}}$  household and  $Y_k$  is the value of the variable  $Y$ , defined in step (2) above, for the  $k^{\text{th}}$  household. As you can see, index  $p$ , the province indicator, takes values ranging from 1 to 11. Eleven distinct province codes appear on the microdata file: one for each of the ten provinces, and a "00" province code assigned to a set of records for reasons of confidentiality. (See Confidentiality of the public-use microdata on page 30.)

**Important:** When estimating variance for a given domain, do not limit yourself to units belonging to the domain. The entire sample should always be used to estimate variance. Units that do not belong to the domain of interest are not considered when computing the point estimate of the total, but do contribute when estimating the variance.

**Rule 2: Approximating CV for estimates of averages or proportions**

An estimated mean or proportion is obtained by the ratio of two estimated totals. For a proportion, the numerator is an estimate that is a sub-set of the denominator, for example the proportion of expenditures for households in Manitoba compared to all Canadian households. The CV of an estimated mean or proportion tends generally to be slightly lower than the corresponding CV of the numerator. The CV of an estimated mean or proportion can thus be approximated with the CV of the numerator and the technique described in rule (1) can be used.

**Rule 3: Approximating CV for estimates of ratios**

Ratio refers to the relationship between any two estimates of totals for which rule (2) does not apply. Approximate CVs for any other types of ratio, may be calculated using the following formula:

$$ACV_R = \sqrt{ACV_N^2 + ACV_D^2}$$

where  $ACV_R$  is the approximate CV of the ratio,  $ACV_N$  is the approximate CV of the numerator of the ratio and  $ACV_D$  is the approximate CV of the denominator of the ratio. The formula will tend to overestimate the CV if the two estimates forming the ratio are positively correlated and underestimate the CV if these two estimates are negatively correlated.

**Rule 4: Approximating CVs for estimates of differences**

The approximate CV of a difference between any two estimates ( $EST_{DIFF} = EST_1 - EST_2$ ) is given by:

$$ACV_{DIFF} = \frac{\sqrt{(EST_1 ACV_1)^2 + (EST_2 ACV_2)^2}}{|EST_{DIFF}|}$$

where  $ACV_1$  is the approximate CV associated with  $EST_1$  and  $ACV_2$  is the approximate CV associated with  $EST_2$ . The formula will tend to overestimate the CV if the two estimates forming the difference are positively correlated and underestimate the CV if these two estimates are negatively correlated.

**Examples**

Detailed calculations of approximate CVs used for estimating totals are initially presented using fictional cases. Then actual cases of estimating totals, averages (or proportions) ratios and differences, based on microdata file data, will be presented so users can check results and ensure that the method used was valid.

**Part 1: Fictional case: details of calculating an approximated CV for estimating a total**

**A) Quantitative variable**

Let us assume we wanted to estimate the total for a (quantitative) expenditure variable X, for households containing at least one person less than 18 years of age. To illustrate this procedure, we will use a fictional sample (see Figure 3) on which we will present calculation details (see Figure 4) for each of the eleven steps described above. As this procedure is applied independently within each province, we shall merely describe calculations for one province.

Let us use the following sample for Ontario:

**Figure 3**  
**Fictional example**

Initial Data					
Identifier	Province	Weight	Household size	Number of children aged 0-17	Variable of Interest X
00001	Ontario	5	3	2	30
00002	Ontario	20	5	3	0
00003	Ontario	25	2	1	20
00004	Ontario	5	4	2	50
00005	Ontario	15	3	0	20
00006	Ontario	10	1	0	10
00007	Ontario	15	4	0	15

In step 1, we define the domain of interest by creating a binary variable equal to 1 for all units belonging to the domain. In the present case, these are households with at least one child between the ages of 0 and 17 years. We then proceed to steps 2 through 9 to estimate variance, which will lead to calculation of the CV. We thus obtain the following results:

**Figure 4**  
**Calculation details for approximating the CV of a total (steps 1 to 9)**

	Step 1	Step 2	Step 4	Step 5	Step 6	Step 7	Step 8	Step 9
Ident.	Binary variable I	Quantitative variable Y (X * I)	Weighted Y (Weight * Y)	Variable K (Weight * size)		Step 6 * size	(Y - step 7)	(Weight - 1) * Weight * (Step 8) <sup>2</sup>
00001	1	30 * 1 = 30	5 * 30 = 150	5 * 3 = 15		3 * 3 = 9	30 - 9 = 21	(4 * 5 * 21 * 21) = 8,820
00002	1	0 * 1 = 0	20 * 0 = 0	20 * 5 = 100		3 * 5 = 15	0 - 15 = -15	(19 * 20 * (-15) * (-15)) = 85,500
00003	1	20 * 1 = 20	25 * 20 = 500	25 * 2 = 50		3 * 2 = 6	20 - 6 = 14	(24 * 25 * 14 * 14) = 117,600
00004	1	50 * 1 = 50	5 * 50 = 250	5 * 4 = 20		3 * 4 = 12	50 - 12 = 38	(4 * 5 * 38 * 38) = 28,880
00005	0	20 * 0 = 0	15 * 0 = 0	15 * 3 = 45		3 * 3 = 9	0 - 9 = -9	(14 * 15 * (-9) * (-9)) = 17,010
00006	0	10 * 0 = 0	10 * 0 = 0	10 * 1 = 10		3 * 1 = 3	0 - 3 = -3	(9 * 10 * (-3) * (-3)) = 810
00007	0	15 * 0 = 0	15 * 0 = 0	15 * 4 = 60		3 * 4 = 12	0 - 12 = -12	(14 * 15 * (-12) * (-12)) = 30,240
			<b>Total: 900</b>	<b>Total: 300</b>	900 / 300 = 3			<b>Total = 288,860</b>

If we wanted to know the CV for Ontario, we would perform the following calculation:

$$CV_{ONT} = 100 * \frac{\sqrt{Variance_{ONT}}}{Estimation_{ONT}} = 100 * \frac{\sqrt{Step\ 9_{ONT}}}{Step\ 4_{ONT}} = 100 * \frac{\sqrt{288860}}{900} = 59.7$$

If we wanted to know the CV for Canada, we would proceed in similar manner, by totalling the results for each province. In other words,

$$CV_{CAN} = 100 * \frac{\sqrt{Variance_{CAN}}}{Estimation_{CAN}}$$

$$= 100 * \frac{\sqrt{\text{Variance}_{NL} + \dots + \text{Variance}_{BC} + \text{Variance}_{PROV00}}}{\text{Estimation}_{NL} + \dots + \text{Estimation}_{BC} + \text{Estimation}_{PROV00}}$$

## B) Qualitative variable (categorical)

In the event a categorical variable is estimated, the steps in calculating the approximate CV will be the same as in the quantitative variable example presented. Instead of a quantitative value for variable of interest X, we would create a dichotomous variable that would be equal to 1 if the household has the features we want to estimate. If not, it would be equal to 0.

To estimate categorical variables, various approaches may be used for defining the domain and the variable of interest, both of which will produce the same results.

Let us assume we want to estimate the number of households consisting of more than one person living in a single-family dwelling. We could proceed in different ways:

- 1) Binary variable I is equal to 1 for all households and variable X is equal to 1 for households consisting of more than one person living in a single-family dwelling.
- 2) Binary variable I is equal to 1 for all households consisting of at least one person and variable X is equal to 1 for all households the members of which live in a single-family dwelling.
- 3) Binary variable I is equal to 1 for all households the members of which live in a single-family dwelling and variable X is equal to 1 for all households made up of more than one person.
- 4) Binary variable I is equal to 1 for all households made up of more than one person living in a single-family dwelling and X is equal to 1 for all households.

Whatever approach is used, the resulting Y variable (step 2) will be equal to 1 if the household possesses all the necessary features (more than one person and living in a single-family dwelling). If not, it will be equal to 0. Results in terms of point estimates and estimates of variance (CV) will thus be the same.

## Part 2: Actual cases based on the microdata file

### Example 1a: Approximation of CV for estimates of totals (quantitative variable)

Let us assume that we have estimated that household furnishings and equipment expenditures for one-person households in Manitoba total \$108,978,360. We have to estimate the approximate CV for this estimate. Users must therefore follow steps (1) to (11) of rule 1.

- 1) Create a binary variable I whose value is 1 if the household is a one-person household and resides in Manitoba, otherwise I equals 0.
- 2) Y is defined for each household as the product of the binary variable I and the 'total household furnishing and equipment expenditures' variable.

Note that the estimate of spending on household furnishings and equipment is obtained by adding the product of variable Y defined in 2) and the final weight of the household.

Figure 5 shows the results of some of the steps in the approximate CV calculation.

**Figure 5**  
**Calculation of ACV**

Step	Total spending on household furnishings and equipment for one-person households in Manitoba
4	108,978,360
5	1,080,391
6	100.87
9	$1.6759 \times 10^{14}$
10	$1.6759 \times 10^{14}$
11	11.88

**Example 1b: Approximation of CV for estimates of totals (qualitative variable)**

Let us assume we now want to estimate the total number of Canadian one-person households, as well as the total number of Canadian households made up of one person living in different types of accommodations.

In this case, variable I is defined as having the value 1 if the household is one-person. If not, it is 0. We must create five Z variables: Z1 with a value of 1 if the type of residence occupied is a "single-family dwelling," and 0 if not; Z2 equals 1 if the type of residence is semi-detached, and 0 if it is not. Z3 equals 1 if the type of residence is a townhouse, and 0 if it is not. Z4 equals 1 if the type of residence is a row house, and 0 if it is not. Finally, Z5 equals 1 if the type of house is "other," and 0 if it is not. Y1 is defined as the product of I and Z1, Y2 as the product of I and Z2, etc.

The estimates obtained are 3,544,346 for the set of one-person households, 1,117,096 for single-family dwellings<sup>3</sup>, 131,821 for semi-detached houses<sup>4</sup>, 183,953 for town houses<sup>5</sup> and 2,111,476 for "other"<sup>6</sup>. We want to calculate the approximate CVs for these estimates.

3. Single family = single detached

4. Semi-detached = double

5. Town houses = row or terrace

6. Other = duplex, apartment, hotel, mobile home, other

Figure 6 shows the results for some steps in the calculation of the approximate CV. The results presented for steps 4 to 9 are the results for Manitoba (presented as an example, for a province, they will be used for comparison in the next example), while those presented for steps 10 and 11 are Canada-wide.

**Figure 6**  
**Calculation of ACV**

Step	Number of one-person households	Number of one-person households living in a single-family dwelling	Number of one-person households living in a semi-detached dwelling	Number of one-person households living in a townhouse	Number of one-person households living in other housing
4	135,175	57,667	2,676	3,242	71,590
5	1,080,391	1,080,391	1,080,391	1,080,391	1,080,391
6	0.13	0.05	0.002	0.003	0.07
9	53,208,085	21,532,096	856,331	1,092,237	26,156,116
10	7,208,220,960	1,943,714,331	254,902,699	326,685,449	4,266,746,906
11	2.40	3.95	12.11	9.83	3.09

**Example 1c: Approximation of CV for estimates of totals used in the calculation of average expenditure**

Let us assume we want to estimate average expenditure on furnishings and household equipment for one-person households in Manitoba. To do so, we would have to estimate the number of one-person households in Manitoba, as well as the total of their expenditure on furnishings and household equipment.

**Figure 7**  
**Calculation of ACV**

Step	Number of one-person households in Manitoba	Total expenditure on furnishings and household equipment for households consisting of one person in Manitoba
4	135,175	108,978,360
5	1,080,391	1,080,391
6	0.13	100.87
9	53,208,085	$1.6759 \times 10^{14}$
10	53,208,085	$1.6759 \times 10^{14}$
11	5.40	11.88

The estimate of the mean would be  $\$108,978,360/135,175 = \$806.2$  How do we determine the CV of this estimate?

Rule (2) should be applied in this case. Thus, the CV of this mean may be approximated with the CV of the numerator, the total expenditure on furnishings and household equipment in Manitoba for one-person households. This CV is 11.88%.

**Example 2: Approximation of CV for estimating ratios**

Let us assume we want to estimate the ratio between the total expenditures on furnishings and household equipment for couples without children households in urban Manitoba and rural Manitoba.

**Figure 8  
Calculation of ACV**

Step	Total expenditure on furnishings and household equipment for households consisting of couple without children and without additional persons in Manitoba (urban)	Total expenditure on furnishings and household equipment for households consisting of couple without children and without additional persons in Manitoba (rural)
4	159,759,581	79,874,673
5	1,080,391	1,080,391
6	147.87	73.93
9	$2.5098 \times 10^{14}$	$1.8193 \times 10^{14}$
10	$2.5098 \times 10^{14}$	$1.8193 \times 10^{14}$
11	9.92	16.89

The estimate of the ratio would be equal to  $\$159,759,581 / \$79,874,673 = 2.0$  (couple without children households in urban Manitoba spend approximately 2 times more on furnishing than those in rural Manitoba). How does the user determine the CV of this estimate?

We have already calculated CVs for each of the two estimates involved in estimating the ratio. We would thus apply rule (3) to obtain the desired CV:

$$CVA_R = \sqrt{CVA_N^2 + CVA_D^2} = \sqrt{9.92^2 + 16.89^2} = 19.59$$

This CV should be identified as “Marginal” (see Guidelines for release) as it is quite high, being between 16.6% and 33.3%.

**Example 3: Approximation of CV for estimating differences**

Let us assume we wanted to estimate the difference between total expenditures on furnishings and household equipment in Alberta and in Manitoba, as well as the CV for this difference.

We would estimate total expenditures on furnishings and household equipment, along with their respective CVs for Manitoba (total = \$741,389,499; CV = 3.82) and for Alberta (total = \$3,682,251,437; CV = 4.49).

Estimation of the difference would thus be \$3,682,251,437 – \$741,389,499 = \$2,940,861,938. Rule (4) can be applied to obtain the desired CV.

$$CVA_{DIFF} = \frac{\sqrt{(EST_1 CVA_1)^2 + (EST_2 CVA_2)^2}}{|EST_{DIFF}|}$$

$$= \frac{\sqrt{(3,682,251,437 * 4.49)^2 + (741,389,499 * 3.82)^2}}{|2,940,861,938|} = 5.70$$

### How to obtain confidence limits

Although coefficients of variation are widely used, a more intuitively meaningful measure of sampling error is the confidence interval of an estimate. A confidence interval constitutes a statement on the level of confidence that the true value for the population lies within a specified range of values. For example a 95% confidence interval can be described as follows.

If sampling of a population is repeated many times, each sample leading to a new confidence interval for an estimate, then in 95% of the samples the interval will cover the true population value.

Using the CV of an estimate, its confidence intervals may be obtained assuming that, under repeated sampling of the population, the various estimates obtained for a characteristic are normally distributed around the true population value. Using this assumption, the chances are about 68 out of 100 that the difference between a sample estimate and the true population value would be less than one standard error, about 95 out of 100 that the difference would be less than two standard errors, and about 99 out 100 that the differences would be less than three standard errors. These different degrees of confidence are referred to as the confidence levels.

Confidence intervals for an estimate, EST, are generally expressed as two numbers, one below the estimate and one above the estimate, as (EST - k, EST + k) where k is determined depending on the level of confidence desired and the sampling error of the estimate.

Confidence intervals for an estimate can be calculated by first determining the ACV of the estimate and then using the following formula to convert to a confidence interval CI:

$$(EST - z \times EST \times ACV / 100, EST + z \times EST \times ACV / 100)$$

where

z = 1 if a 68% confidence interval is desired,

$z = 1.6$  if a 90% confidence interval is desired,  
 $z = 2$  if a 95% confidence interval is desired,  
 $z = 3$  if a 99% confidence interval is desired.

**Note:** Release guidelines, which apply to the estimate, also apply to the confidence interval. For example, if the estimate is not releasable, then the confidence interval is not releasable either.

#### Example 4

A 95% confidence interval for the estimated mean of spending on household furnishings and equipment for one-person households in Manitoba would be calculated as follows:

$$\text{EST} = \$806.2$$

$$z = 2$$

$$\text{ACV} = 11.88$$

$$\text{CI} = (806.2 - 2 \times 806.2 \times 11.88/100; 806.2 + 2 \times 806.2 \times 11.88/100) = (\$614.6, \$997.8)$$

#### How to do a Z-test

Coefficients of variation may also be used to perform hypothesis testing, a procedure for distinguishing between population parameters using sample estimates. The sample estimates can be totals, averages, ratios, etc. Tests may be performed at various levels of significance, where a level of significance is the probability of concluding that the characteristics are different when, in fact, they are identical.

Let  $\text{EST}_1$  and  $\text{EST}_2$  be sample estimates for 2 characteristics of interest. Let the approximate CV of the difference  $\text{EST}_1 - \text{EST}_2$  be  $\text{ACV}_{\text{DIFF}}$ .

If  $z = 100 / \text{ACV}_{\text{DIFF}}$  is less than 2, then no conclusion about the difference between the characteristics is justified at the 5% level of significance. If however, this ratio is larger than 2, the observed difference is significant at the 5% level.

#### Example 5

Let us suppose we wish to test, at the 5% level of significance, the hypothesis that there is no difference between the total of spending on furnishings and equipment in Alberta and the same total in Manitoba. From example 3, the approximate CV of the difference between these two estimates was found to be 5.70 and  $z = 17.54$ . Since this value is greater than 2, it must be concluded that there is significant difference between the two estimates at the 0.05 level of significance.

## Confidentiality of the public-use microdata

Microdata files for public use differ in many ways from the master file of the survey held by Statistics Canada. These variations are due to measures taken to preserve the anonymity of respondents to the survey.

The confidentiality of this file is ensured mainly by reducing information, i.e., deleting variables or suppressing or collapsing some of their detail.

### To protect confidentiality

- All explicitly identifying information, such as identification numbers, was removed from the file. (Names and addresses are not data captured).
- 205 records had their *province codes* set to 0 due to special characteristics (e.g., exceedingly high or low expenditure values). These records were reweighted.
- Other records were also reweighted for confidentiality reasons.
- There was *top-coding* and *collapsing* of code sets for non-spending variables.
- Income values at the household, reference person and spouse of reference person levels were *rounded* in the following manner:
  - For income values between \$1 and \$9,999: round to the nearest \$100
  - For income values between \$10,000 and \$99,999: round to the nearest \$1,000
  - For income values between \$100,000 and \$999,999: round to the nearest \$10,000
  - For income values between \$1,000,000 and \$9,999,999: round to the nearest \$100,000
  - For income values between \$10,000,000 and \$99,999,999: round to the nearest \$1,000,000 (there are no such values on the 2006 file).
- The variables “Purchase price of dwelling” and “Selling price of dwelling” were also rounded.

## **Appendices—See accompanying Excel file**

### **Appendix A Frequency counts**

### **Appendix B Averages, aggregates, minimum and maximum values**

Part 1 of 2 – Suppressed PUMF file

Part 2 of 2 - Unsuppressed survey file

### **Appendix C Inclusion of spending variables in past microdata files**

### **Appendix D Comparison of variables from the 2005 and 2006 Survey of Household Spending**

### **Appendix E Coefficients of variation for published data from the 2006 SHS**

Part 1 of 3 - Average expenditure per household, Canada and provinces

Part 2 of 3 - Median expenditure per household reporting, Canada and provinces

Part 3 of 3 - Dwelling characteristics and household equipment, Canada and  
provinces