

## TABLE OF CONTENTS

	Page
1. Introduction .....	7
2. Background .....	8
3. Objectives .....	8
4. Survey Content .....	9
4.1 Selection Criteria .....	9
4.2 1996-97 Changes to Existing Content .....	10
4.3 1994-95 Data Feedback and Follow-up Questions .....	12
4.4 New Content for 1996-97 .....	13
5. Sample Design .....	15
5.1 1996-97 Sample Design .....	15
5.1.1 Sample Design for the Core Household Component .....	15
5.1.2 1996-97 Sample Design for Provincial Supplements to the Household Component .....	19
5.1.3 Overview of Random Digit Dialing Sampling .....	19
5.1.4 Sample Size Allocation by Health Area .....	20
5.1.5 Stratification and Sample Allocation .....	21
5.2 1994-95 Sample Design .....	22
5.2.1 Sample Design for the Household Component .....	22
5.2.2 Sample Allocation .....	22
5.2.3 The Rejective Approach .....	23
5.2.4 Sample Selection .....	24
5.2.5 Sample Design in Québec .....	25
6. Data Collection .....	27
6.1 Questionnaire Design and Data Collection Method .....	27
6.2 Tests .....	27
6.3 Interviewing .....	28
6.4 Supervision and Control .....	28
6.5 Non-response to the NPHS .....	29
6.6 Non-response Follow-up .....	29
6.7 Tracing .....	30
7. Data Processing .....	31
7.1 Editing .....	31
7.2 Coding .....	31
7.3 Creation of Derived and Grouped Variables .....	31
7.4 Weighting .....	31
7.5 Suppression of Confidential Information .....	32

8.	Data Quality .....	33
8.1	Response Rates .....	33
8.1.1	Core Cross-sectional Response Rates .....	34
8.1.2	RDD Cross-sectional Response Rates (by province and total) .....	35
8.1.3	Overall Cross-sectional Response Rates .....	37
8.2	Survey Errors .....	39
9.	Guidelines For Tabulation, Analysis And Release .....	41
9.1	Rounding Guidelines .....	41
9.2	Sample Weighting Guidelines for Tabulation .....	42
9.2.1	Definitions of Types of Estimates: Categorical vs. Quantitative .....	42
9.2.2	Tabulation of Categorical Estimates .....	43
9.2.3	Tabulation of Quantitative Estimates .....	43
9.3	Guidelines for Statistical Analysis .....	44
9.4	Release Guidelines .....	45
10.	Approximate Sampling Variability Tables .....	47
10.1	How to Use the C.V. Tables for Categorical Estimates .....	53
10.2	Examples of Using the C.V. Tables for Categorical Estimates .....	55
10.3	How to Use the C.V. Tables to Obtain Confidence Limits .....	59
10.4	Example of Using the C.V. Tables to Obtain Confidence Limits .....	60
10.5	How to Use the C.V. Tables to do a Z-test .....	60
10.6	Example of Using the C.V. Tables to do a Z-test .....	60
10.7	Exact Variances/Coefficients of Variation .....	61
10.8	Release Cut-offs for the NPHS .....	62
11.	Weighting .....	69
11.1	Cross-sectional Weighting for the 1996-97 NPHS—Core Household Sample .....	70
11.1.1	Stripped Weights .....	70
11.1.2	Weight Adjustments for Household Members .....	70
11.1.3	Weight Adjustments for Selected Members .....	74
11.2	Cross-sectional Weighting for the 1996-97 NPHS—Provinces with RDD Supplemental Samples .....	76
11.2.1	RDD Basic Weights .....	77
11.2.2	Further Weight Adjustments to the Basic Weights .....	78
11.2.3	Further Weight Adjustments for Household Members .....	80
11.2.4	Further Weight Adjustments for Selected Members .....	82
11.3	1994-95-based Weighting Procedures for the Provinces Outside Québec .....	86
11.3.1	LFS Basic Weights .....	86
11.3.2	Further Weight Adjustments to the Basic Weights .....	86
11.3.3	Further Weight Adjustments for Selected Members .....	89
11.4	1994-95-based Weighting Procedures for Québec .....	91
11.4.1	ESS Weights .....	91
11.4.2	NPHS Basic Dwelling Weights .....	92

11.4.3	Further Weight Adjustments to the Basic Weights	93
11.4.4	Further Weight Adjustments for Selected Members	94
12.	File Usage	97
12.1	Use of Weights	97
12.1.1	Cross-sectional Weight - General File WT56	97
12.1.2	Cross-sectional Weight - Health File WT66	97
12.1.3	Cross-sectional Weight - Health File WT66_N	98
12.2	Variable Naming Convention	99
12.2.1	Variable Name Component Structure	99
12.2.2	Positions 1-2: Variable / Questionnaire Section Name	100
12.2.3	Position 3: Survey Type	101
12.2.4	Position 4: Year / Cycle Variable	101
12.2.5	Position 5: Variable Type	102
12.2.6	Positions 6-8: Variable Name	102
12.3	Remote Access of Master Files	103

Appendix A: Questionnaire

Appendix B: Record Layout - General Microdata File

Appendix C: Record Layout - Health Microdata File

Appendix D: Data Dictionary Alphabetical Index - General Microdata File  
 Data Dictionary Topical Index - General Microdata File  
 Data Dictionary - General Microdata File

Appendix E: Data Dictionary Alphabetical Index - Health Microdata File  
 Data Dictionary Topical Index - Health Microdata File  
 Data Dictionary - Health Microdata File

Appendix F: Derived and Grouped Variables

Appendix G: C.V. Tables List - General Microdata File  
 C.V. Tables - Canada by Agegroup - General Microdata File  
 C.V. Tables by Province and Canada Total - General Microdata File  
 C.V. Tables by Health Area for Ont., Man. and Alta.- General Microdata File

Appendix H: C.V. Tables List - Health Microdata File  
 C.V. Tables - Canada by Agegroup - Health Microdata File  
 C.V. Tables by Province and Canada Total - Health Microdata File  
 C.V. Tables by Health Area for Ont., Man. and Alta. - Health Microdata File

## **1. Introduction**

The National Population Health Survey (NPHS) Program is designed to collect information related to the health of the Canadian population. The first cycle of data collection began in 1994, and will continue every second year thereafter. The survey program will collect not only cross-sectional information, but also data from a panel of individuals at two-year intervals. It is composed of three component parts: the survey of households; the survey of institutions and the survey of the North.

The household component includes household residents in all provinces, with the principal exclusion of populations on Indian Reserves, Canadian Forces Bases and some remote areas in Québec and Ontario. The institutional component includes long-term residents (expected to stay longer than six months) in health care facilities with four or more beds in all provinces with the principal exclusion of the Yukon and the Northwest Territories. The northern component includes household residents in both the Yukon and the Northwest Territories with the principal exclusion of populations on Indian Reserves, Canadian Forces Bases and some of the most northerly remote areas of the Territories.

The household component of the National Population Health Survey Program has completed two release cycles, 1994-95 and 1996-97. This document has been produced to facilitate the manipulation of the two 1996-97 cross-sectional microdata files containing the survey results. These files are described in more detail in Chapter 4 and the appendices.

Any questions about the data sets or their use should be directed to:

For technical/general data support call:

Electronic Products Help Line 1-800-949-9491

For custom tabulations/general data support call:

Client Custom Services, Health Statistics Division 1-613-951-1746

For remote access support call:

Sylvie Alary, NPHS 1-613-951-1653  
Internet: [nphs@statcan.ca](mailto:nphs@statcan.ca)

For survey content support call:

Bryan Lafrance, NPHS 1-613-951-3285  
Fax: 1-613-951-4198

## **2. Background**

In the fall of 1991, the National Health Information Council (NHIC) recommended that an ongoing national survey of population health be conducted. This recommendation was based on consideration of the economic and fiscal pressures on the health care system and the commensurate requirement for information with which to improve the health status of the population in Canada. Existing sources of health data were unable to provide a complete picture of the health status of the population and the myriad factors that have an impact on health.

Commencing in April 1992, Statistics Canada received funding for development of a National Population Health Survey. The survey was designed to be flexible and to produce valid, reliable and timely data. Also, it was to be responsive to changing requirements, interests and policies.

## **3. Objectives**

The objectives of the NPHS are to:

- C aid in the development of public policy by providing measures of the level, trend and distribution of the health status of the population;
- C provide data for analytic studies that will assist in understanding the determinants of health;
- C collect data on the economic, social, demographic, occupational and environmental correlates of health;
- C increase the understanding of the relationship between health status and health care utilization, including alternative as well as traditional services;
- C provide information on a panel of people who will be followed over time to reflect the dynamic process of health and illness;
- C provide the provinces and territories and other clients with a health survey capacity that will permit supplementation of content or sample;
- C allow the possibility of linking survey data to routinely collected administrative data such as vital statistics, environmental measures, community variables, and health services utilization.

## **4. Survey Content**

The objectives described in Chapter 3 provided only a broad direction for the NPHS, particularly concerning the type of information to be collected. The first section discusses the general criteria used for the selection of survey content and gives a broad summary of sections and changes. This is followed by a section describing detailed changes to existing content for 1996-97. The next section focuses on the 1994-95 variables that were fed back and used in cycle 2. The last section details new content for 1996-97.

### **4.1 Selection Criteria**

Survey content was selected according to the following criteria:

- 1) Information should relate to, and help monitor, the health goals and objectives of the provinces and territories. Where health goals have not been established, for example, at the national level, policies and programs could be considered in the selection of survey content.
- 2) The information should not duplicate data available from other sources.
- 3) With a view to increasing the understanding of health and its determinants, information collected should provide new knowledge in areas that have not been adequately studied.
- 4) The survey should focus on behaviours or conditions amenable to prevention, treatment, or intervention.
- 5) The survey should collect information about conditions that impose the greatest burden, in terms of suffering or cost, on affected individuals, the general population, or the health care system.
- 6) The survey should collect information on factors related to good health, not just those related to illness.

In each household, some limited information was collected from all household members (general component) and one person in each household was randomly selected for a more in-depth interview (health component). Reflecting the above criteria, the questionnaire included sections on health status, use of health services, risk factors and demographic and socio-economic status. For example, health status was measured through questions on self-perception of health, functional ability, chronic conditions, and activity restriction. The use of health services was measured through questions on visits to health care providers, hospital care and drug use. Behavioural risk factors included smoking, alcohol use and

physical activity. In addition, a special focus of this second cycle of the survey was access to services. Questions were asked on preventive tests and examinations which probed for frequency, reasons for use or non-use and barriers encountered. Demographic and socio-economic information included age, sex, education, ethnicity, household income and labour force status.

In order to reduce the collection costs, the 1996-97 Health Promotion Survey (HPS) questions were integrated into the NPHS rather than being done as a separate supplemental survey as in 1994-95. These questions were asked of the core sample and most of the Random Digit Dialing (RDD) sample (see Chapter 5 - Sample Design for an explanation of core and RDD sample). Additional detailed questions covered the following topics: actions taken to improve health, preferred weight, breast-feeding, nutrition, HIV/AIDS, smoking during pregnancy, opinions on smoking, opinions on alcohol, social support, sexual health and road safety.

As part of the buy-in agreement with the Alberta Ministry of Health, supplementary questions were asked of all Alberta respondents (both core and RDD) on topics such as sun exposure, sources of health information, violence and personal safety and several other topics. In addition, some questions that were part of the Health Promotion Survey were not asked of the Alberta RDD sample but were asked of the core Alberta respondents.

As part of buy-in agreements, in both the Manitoba and Alberta RDD samples, a child was selected (where present) in addition to the usual adult. Questions on child health services were answered for these children on their behalf.

A list of the questions is provided in Appendix A.

## **4.2 1996-97 Changes to Existing Content**

The 1996-97 NPHS was collected mainly by telephone whereas the previous cycle was collected primarily by personal interview. For that reason, an effort has been made to make the questions easier to read (for the interviewers) and easier to understand (for the respondents) by using more colloquial wording where possible. Without revamping sections, screens have been redesigned, response categories have been shortened and frequent responses in "other specify" have become itemized categories.

**General Component** - changes by section:

**Health Care Utilization:**

UTIL-Q3 had frequent responses to “other specify” from 1994-95 added as new categories.

UTIL-Q4 had the self-help category removed as an alternative medicine and asked as a separate question, UTIL-Q4A.

UTIL-Q6 had the concept changed. The question has been refined to get the view of the person for whom the proxy is reporting and the proxy's view if the person is younger than 18.

UTIL-Q7 and UTIL-Q10 were converted to list type questions based on the most common responses from 1994-95.

**Restriction of Activities:**

RESTR-Q3 and RESTR-Q4 were combined into one question with a computer fill providing the appropriate wording.

RESTR-Q6 was clarified and went from a mark-all to each type of help being asked separately.

**Chronic Conditions:**

CHR-Q1 was a mark-all question and has been changed to a series of individual prompts for each condition. Bowel disorders and thyroid conditions were added to the list. The sub-questions on type of cancer and asthma wheezing were dropped.

**Socio-demographic Characteristics:**

SOCIO-Q1 had Holland added to the category Netherlands.

SOCIO-Q3 had its concept changed slightly by asking when the respondent first came to Canada to live instead of year of immigration.

SOCIO-Q7 had the race categories reordered and refined to reflect those of the 1996 Census.

**Education:**

EDUC-Q5 and EDUC-Q6 in 1994-95 were moved to the front of the section to become EDUC-Q1 and EDUC-Q2 to allow for the feedback of previous data. As a result, EDUC-Q1 to EDUC-Q4 have become EDUC-Q4 to EDUC-Q7 respectively.

**Labour Force:**

This section was revamped to take advantage of the computer-assisted interviewing (CAI) software capabilities and some results from 1994-95. The section is restricted to persons 15 to 75 years old. The number of jobs collected was reduced from 6 to 3. The main activity question, LFS-Q1, was dropped in favour of having the application decide which was the main job. A gap in employment was redefined as 28 days instead of 7 days. LFS-Q17A and LFS-Q17B had four new categories added: retired, resigned, looking for work, and disabled/recovering from illness.



**Health Component** - changes by section:

The sections were rearranged to enhance the presentation of the questions especially with the integration of the Health Promotion Survey (HPS) and the Alberta supplement.

**Height and Weight:**

The minimum and maximum weights for adults allowed in CAI were changed from 1 and 999 to a range of 18-575 based on 1994-95 results.

**Access to Services:**

Questions on blood pressure, mammography, and pap smear testing were incorporated into this section.

**Physical Activities:**

In PA-Q1, the skating category was refined to ice skating. Inline or roller skating would now be included in other (specify). The category tai-chi was replaced by basketball.

**Injuries:**

IN-Q7 was only asked if the respondent was currently working.

**Drug Use:**

DRG-Q1 was a mark-all question and has been changed to a series of individual prompts for each drug category. The categories for hormone use and birth control pills were reversed. A new category for thyroid medication was added.

**Alcohol:**

ALC-Q2 had the order of the response categories reversed.

ALC-Q3 had the response categories converted to a list of frequencies, e.g., once a month, instead of collecting an actual number.

**Social Support:**

SUP-Q7 had the category "don't have any" changed to "don't have any or all live with you".

**Health number and Administration:**

The longitudinal respondent was asked if his/her health number had changed since the last interview. If yes, we collected it otherwise the question was skipped.

### **4.3 1994-95 Data Feedback and Follow-up Questions**

In order to reduce respondent burden, questions to which we already knew the answer that would not change over time (e.g., country of birth) were not asked again. For variables that could change over time but only if certain actions had occurred (e.g., level of education), updating was only done if appropriate.

**Restriction of Activities:**

Whether or not the respondent had a disability in 1994 was used. If the status changed, we probed for an explanation of that change.

**Chronic Conditions:**

For all respondents, selected chronic conditions (asthma, arthritis, high blood pressure, migraine headaches, diabetes, epilepsy, stomach or intestinal ulcers and the effects of a stroke) were fed back in an attempt to help explain change. If a newly acquired condition, the date of onset for the condition was acquired.

**Socio-demographic Characteristics:**

For all respondents, flags indicating that country of birth and ethnic origin had been collected were re-input. Since the response categories to race were changed, this variable was re-collected. Language first learned and still spoken was asked again because it can change over time.

**Education:**

For all respondents, a flag indicating the highest level of education was re-input. Screening questions determined if the respondent was currently attending a learning institution between cycles. If so, educational attainment was collected anew.

**Labour Force:**

For all respondents, the employer name, type of industry and duties for the main job in 1994-95 were fed back. If the respondent indicated that they worked in the previous year, they were asked to confirm the employer name.

**Smoking:**

For the longitudinal respondent, type of smoker was fed back. If it was a non-proxy interview in 1994-95, reasons for change were probed.

#### **4.4 New Content for 1996-97**

**General Component:**

**Health Care Utilization:**

Two new questions were included on health care services received in the United States.

**Chronic Conditions:**

For certain chronic conditions (asthma, arthritis/rheumatism, high blood pressure, migraine headaches and diabetes) follow-up questions on whether any medication or treatment was being taken plus the type of treatment/medication taken was collected.

**Health Component:**

**Access to Services:**

In 1996-97 the core focus was on access to services. For various health services, questions on the frequency of use, reasons for use, barriers encountered, and reasons for non-use or less frequent use than recommended by the Canadian Medical Association were probed. These services included: blood pressure, pap smear test, mammography, physical check-up, flu shots, dental visits and eye examinations.

**Repetitive Strain:**

Questions on repetitive strain were new and asked immediately before the injury section so as to separate out this type of strain injury.

**Drug Use:**

If hormone use was reported, questions were asked about the type of hormones taken and the year in which the treatment began.

**Alcohol Dependence:**

This is part of the Kessler and Mroczek series on the Composite International Diagnostic Interview (CIDI) 12-month Short Screening Scales.

**Mental Health:**

A follow-up question was added for those who reported seeing a doctor in the past 12 months. It was a pick-list of whom they saw or to whom they talked.

**Administration:**

The contact information for the longitudinal respondent was collected as part of the health component. The telephone number at work was a new variable collected.

## **5. Sample Design**

The target population of the NPHS includes household residents in all provinces, with the principal exclusion of populations on Indian Reserves, Canadian Forces Bases and some remote areas in Québec and Ontario.

### **5.1 1996-97 Sample Design**

#### **5.1.1 Sample Design for the Core Household Component**

In the first cycle of the NPHS, the sample was created by first selecting households and then within each household choosing one member to be the longitudinal respondent. For the second cycle, the distinction is made between the sample selected for longitudinal purposes and the sample selected for cross-sectional purposes.

The longitudinal sample for 1996-97 consists of all longitudinal respondents chosen in cycle 1 who had completed at least the general component of the questionnaire in 1994-95. This included 2,022 persons who were under the age of 12 in cycle 1 (previously interviewed as part of the 1994-95 National Longitudinal Survey of Children and Youth (NLSCY) who were included in the NPHS sample for 1996-97). Units selected in 1994-95 as part of supplemental buy-in samples were excluded. Only the longitudinal respondent was traced using contact information collected in 1994-95; no attempt was made to follow all household members over time. For cross-sectional purposes, all household members currently living with the longitudinal respondent were interviewed. The selected longitudinal respondent's data will be used for longitudinal purposes and cross-sectional purposes.

The core sample selected in 1994-95 was not increased for 1996-97 cross-sectional estimates. However, in three provinces, there are supplementary buy-in samples for cross-sectional purposes (see 5.1.2). In the provinces without supplementary samples, the sample size is slightly diminished from 1994-95 due to attrition of the sample, deaths, out-of-scope cases, untraceable cases and non-respondents. This decrease in sample size is relatively small, and should not lead to large increases in the variance of estimates. Analysis of the 1996-97 sample before data collection showed that the sample was not biased in terms of representativity by province, age or sex.

**Longitudinal Core Sample Size by Current Province**  
Number of Persons

Province	1994-95 Core Sample <sup>1</sup>	1996-97 Core Sample	% of 1994-95 Core Sample Followed
Newfoundland	1,212	1,082	89.3
Prince Edward Island	1,184	1,037	87.6
Nova Scotia	1,271	1,085	85.4
New Brunswick	1,277	1,125	88.1
Québec	3,430	3,000	87.5
Ontario	5,335	4,307	80.7
Manitoba	1,346	1,205	89.5
Saskatchewan	1,320	1,168	88.5
Alberta	1,697	1,544	91.0
British Columbia	2,023	1,723	85.2
Total	20,095	17,276	86.0

Since no new longitudinal units were selected in cycle 2, the population covered by the longitudinal sample in 1996-97 is 2 years old and above. That is to say, no selected person is of age 0 or 1 in cycle 2, and therefore, no direct estimates of variables contained on the core health component of the questionnaire will be possible for these ages. Another implication of not selecting any new longitudinal units in cycle 2 is that people who immigrated to Canada in the last two years (since the last time longitudinal units were selected in cycle 1) are not represented by the selected persons forming the longitudinal sample. However, estimates of general component variables are available for the entire 1996-97 cross-sectional population, including the 0 to 1-year-old group. Recent immigrants who joined households that participated in the survey in 1994-95 are also represented by the core sample. It should be noted that post-stratification to 1996-97 population totals (see Chapter 11 - Weighting) implicitly treats the non-covered immigrant population as if it had the same characteristics as the rest of the population.

---

<sup>1</sup> Effective sample, excluding dwellings that were either non-eligible, vacant, under construction or rejected by the rejective method.

**Cross-sectional Core Sample Size by Current Province  
(excluding provincial supplements)  
Number of Persons**

Province	General Component	Health Component
Newfoundland	3,017	1,082
Prince Edward Island	2,752	1,037
Nova Scotia	2,775	1,085
New Brunswick	2,888	1,125
Québec	7,838	3,000
Ontario	10,899	4,307
Manitoba	3,045	1,205
Saskatchewan	2,785	1,168
Alberta	4,164	1,544
British Columbia	4,276	1,723
Total	44,439	17,276

The target population in 1994-95 was household residents in all provinces excluding persons living on Indian Reserves, Canadian Forces Bases and remote areas in Ontario and Québec. Cross-sectionally, the core target population includes all household members currently living with longitudinal respondents. The target population for the provincial supplements is described in section 5.1.2.

To identify whether persons were part of the target population or not, the following rules were developed for data collection purposes.

**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

<b>Status of Longitudinal Respondent</b>	<b>Action Taken</b>
Dead	For longitudinal respondents identified to be deceased at the time of the cycle 2 interview, the death was confirmed against the Canadian Mortality Database. Longitudinal respondents who have died are part of the longitudinal file, but do not appear on cross-sectional files.
Moved into an institution	Longitudinal respondents who moved from a private household to a health care institution were interviewed by the Institutional component of the NPHS. Their data appear on the longitudinal household file but they are not part of the cross-sectional file.
Moved to Northwest Territories or Yukon	In cases where the longitudinal respondent moved, attempts were made to collect their new location. When possible, the longitudinal respondent was interviewed using the household questionnaire. They are part of the longitudinal file, but do not appear on the cross-sectional files.
Moved to Indian Reserve or Canadian Forces Base	If a longitudinal respondent was traced to an Indian Reserve, or Canadian Forces Base, attempts were made to interview the selected person and the other members of their household using the household questionnaire. They are part of the longitudinal file but do not appear on the cross-sectional file.
Moved out of Canada - temporarily	In cases where the longitudinal respondent moved out of Canada for a period of time but is expected to return to living in Canada, attempts were made to interview the longitudinal respondent. Again, these units form part of the longitudinal file, but do not appear on the cross-sectional file.
Moved out of Canada - permanently	If the longitudinal respondent moved out of Canada and is not expected to return, their new location was collected and may be used to follow-up in future waves. These persons were not interviewed and do not appear on the longitudinal or cross-sectional files.

### **5.1.2 1996-97 Sample Design for Provincial Supplements to the Household Component**

In three provinces, Alberta, Ontario and Manitoba, the provincial governments provided extra funds so that a larger sample of dwellings could be selected. The purpose of this buy-in was to get sufficient sample size to provide reliable cross-sectional estimates at sub-provincial (health area) levels. The buy-in sample is combined with the core sample to produce one large cross-sectional file of data.

All of the interviews were done by telephone. A sample of telephone numbers was selected from the Statistics Canada Random Digit Dialing (RDD) System. This means that people without telephones in their homes had no chance of being selected. In addition, people living on Indian Reserves, Canadian Forces Bases or in institutions or collective dwellings were not eligible to be interviewed. Differences in the covered population are adjusted for in the weighting procedures. The sample size per health area was based upon the funding available and the requirements of the provinces to obtain reliable estimates by health area.

As in the core survey, a general component of the questionnaire was administered to all of the members of a responding household. In all three of the provinces, one member of the household aged 12 and over was selected to answer the health component. In Alberta and Manitoba, a child aged 0-11 was also selected (where possible) and administered the same questionnaire as the core selected children.

### **5.1.3 Overview of Random Digit Dialing Sampling**

The Statistics Canada Random Digit Dialing sampling system uses a method called the Elimination of Non-Working Banks (ENWB) which works in the following manner. A bank (the area code plus the first five digits of a telephone number) is considered to be eligible for sampling if it contains at least one residential telephone number as determined from a residential telephone billings file obtained from the telephone companies. Eligible banks are grouped together to form strata.

Within a stratum, a bank is randomly chosen and a number between 00 and 99 is randomly generated. Combining the bank with this random number produces a seven-digit telephone number. A historical file ensures that this number has not been used by another RDD survey in the recent past. This process repeats until the desired number of telephone numbers within the stratum has been generated. Often the generated number corresponds to a telephone number that is not in service. This means that to meet sample size targets, many extra numbers have to be generated. The overall "hit rate" is usually around 50%, meaning that about one-half of the generated numbers actually correspond to households. Of course this varies from



area to area.

#### **5.1.4 Sample Size Allocation by Health Area**

Sample sizes were determined in consultation with the sponsoring ministries of health based upon two criteria, the funding available and the precision of the estimates required. Health areas were defined by the provincial ministries of health using 1991 Census geography.

In Ontario it was decided that each of the 23 health areas required a total of 1,200 health component respondents who agreed to share their data with the Ministry of Health. Exceptions to this rule occurred in the health areas covering Toronto and Ottawa. In these health areas the required number of respondents was 3,000 and 2,000 respectively.

In Manitoba, the number of RDD respondents was set so that the number of core and RDD respondents combined met a certain threshold. In the health area containing Winnipeg, the total of the two components was to reach at least 1,428 respondents who would answer the health component. In the two northern health areas, the total was set to 603. In the other eight areas the total was set to 1,000. The number of respondents expected from the core was about 30 in the north, between 30-110 in the non-Winnipeg southern areas and about 325 in the Winnipeg health area. The RDD component would be responsible for the rest of the respondents.

In Alberta, the required number of respondents was calculated using the following mathematical procedure. In each of the 17 health areas, a required precision of the estimates was defined. The sample size had to be such to allow this precision to be met for a characteristic that occurred in 15% of the population at a sub-health area demographic level. The demographic level was set at either the sex level, the age group level (ages 0-11, 12-24, 25-44, 45-64 and 65+), or the age group/sex level (a crossing of the two previous groupings).

For example, in the Calgary health area there is a requirement that a characteristic found in 15% of the population have a coefficient of variation (C.V.) of no greater than 25% for estimates at the age/sex level. In the northern health areas the requirement was also a C.V. of 25% but this time this precision only had to be met at the sex level. Design effects and the expected number of respondents in each demographic group within the health area were used to determine the required number of households. The number of households necessary ranged from 230 to 2,450 per health area depending upon the requirements.

The table below shows the provincial totals (numbers of respondents and telephone numbers are approximate)

**RDD Supplemental Sample Sizes by Province**

Province	Number of Health Areas in Microdata Files	Number of Health Areas	Targeted Response Rate	Expected Number of Respondents	Number of Generated Telephone Numbers
Ontario	23/16	23	75%	32,000	89,500
Manitoba	5	11	80%	10,300	29,000
Alberta	5	17	80%	13,000	30,300

**5.1.5 Stratification and Sample Allocation**

Telephone exchanges (the first three digits of the telephone number) were used as stratification units and strata were composed of groups of exchanges. By using postal codes found on telephone billing files, a geographic location was determined for the telephone number.

Primary strata were composed of exchanges that covered the same health areas. A rule based upon the percentage of numbers in the exchange falling into each health area was established to determine how many regions an exchange "significantly" covered. Most covered only one, but some covered two, or in a few cases, three health areas. A sampling rate for each region was determined based upon the required number of sampled telephone numbers, the hit rate and the estimated number of telephone numbers in the health area.

In strata that covered only one health area significantly, the stratum was sampled at the rate determined for that health area. In strata that covered more than one health area significantly, the higher sampling rate among the significantly covered health areas was used. Primary strata had to be large enough so that a sample of at least forty households could be expected. If the expected size was smaller than this, the stratum was collapsed with one that covered similar health areas. In larger primary strata, geographic sub-stratification took place. The expected sample sizes from the final strata ranged from 40 to about 500 households.

## **5.2 1994-95 Sample Design**

The redesigned Labour Force Survey (LFS) was used as the basis for the design in all provinces except Québec where the NPHS selected a sample from households already being used by Santé Québec for the 1992-93 *Enquête sociale et de santé* (ESS).

### **5.2.1 Sample Design for the Household Component**

Three factors shaped the design of the household component sample:

- C the targeted national and provincial/territorial sample sizes;
- C the decision to select one member per household to make up the longitudinal panel;
- C the choice of the redesigned LFS as a vehicle for selecting the sample.

These three factors resulted, respectively, in the allocation of the sample, the application of a technique (the "rejective method," described later) to improve the sample's representativeness, and the selection of provincial samples outside Québec.

### **5.2.2 Sample Allocation**

The NPHS was budgeted for a sample size of 19,600 households. It was further agreed among national and provincial representatives that each province needed a minimum of 1,200 households. Subject to this restriction the provincial sample sizes were obtained by using a well-known allocation scheme that balances the reliability requirements at national and regional levels (Kish, 1988)<sup>1</sup>. According to this scheme the sample was allocated proportionally to  $\%(0.804W_h^2 + 1/12^2)$ , where  $W_h$  is the 1991 Census proportion of households in province/territory  $h$ ,  $h=1, \dots, 12$ . This allocation determined the base sample size for each province. Four provinces chose to increase their allotted sample size through the buy-in of additional units.

Within provinces the sample was initially distributed proportionally to the population size. The provincial buy-in samples and the use of a rejective method, described below, affected the sub-provincial allocations. Ontario and Manitoba's buy-in samples imposed minimum requirements by health areas, while New Brunswick and British Columbia paid for additional sample coverage of certain areas only. In B.C. most of the buy-in requirement was met using telephone interviews from a Random Digit Dialing (RDD) sample of telephone numbers. In applying the rejective method, sample sizes were inflated by the number of households expected to be screened out of the sample.

### **5.2.3 The Rejective Approach**

The survey content primarily focused on one member in each sample household who was chosen at random to become the longitudinal panel respondent. Without the use of the rejective method, the panel would underrepresent persons coming from large households, typically parents and children, since they had less chance of being chosen and overrepresent persons coming from small households, often single people or the elderly.

Thus, a rejective approach was adopted to increase the representation of parents and youths in the panel. A portion of the sample was pre-identified for screening. After their member roster was completed, screened households that had no member under 25 years of age were eligible for rejection (EFR) and dropped out of the survey. In order to maintain the required sample sizes, the number of households visited in each province was increased by the anticipated number of households screened out in this way.

The rejective method with an under-25-year-old rule was adopted as it performed better than other rejection rules considered. For cost and operational reasons the percentages of screened households was usually limited to 25-30% in Ontario, 37.5-40% in urban areas elsewhere and 25-30% in rural areas. As apartment strata had a high concentration of small households, their sample sizes were reduced instead of applying a rejective method. The rejective approach was also not applied in remote regions because of the high contact costs there, and its use was limited in areas where sample buy-in demands were substantial.

#### **5.2.4 Sample Selection**

The sample design considered for the household component of the NPHS was a stratified two-stage design. In the first stage homogeneous strata were formed and independent samples of clusters were drawn from each stratum. In the second stage dwelling lists were prepared for each cluster and dwellings, or households, were selected from the lists.

In all provinces except Québec the NPHS used the multi-purpose sampling methodology developed for the redesign of the LFS. That methodology provided general household surveys with clustered samples of dwellings, thus making the design very cost effective for the listing and collection of data.

The basic LFS design is a multi-stage stratified sample of dwellings selected within clusters. Each province is divided into three types of areas (Major Urban Centres, Urban Towns and Rural Areas) from which separate geographic and/or socio-economic strata are formed. In most strata, six clusters, usually Census Enumeration Areas (EAs), were selected with Probability Proportional to Size (PPS). In a few cases where the population density was low an additional stage was added by first selecting two or three large Primary Sampling Units, dividing them into clusters, and drawing a sample of six clusters from each. The number six was used throughout the sample design to allow a one-sixth rotation of the sample every month for the LFS.

The sample of dwellings is obtained after listing operations in sample clusters were completed. As sampling rates were predetermined there were often differences between anticipated and obtained sample counts. Excessive sample yields were corrected by dropping a portion of the originally selected units. This was usually done at aggregated levels and was called sample stabilization. Note also that sample sizes were inflated to represent dwellings rather than households, as approximately 15% of the dwellings were expected to be vacant or otherwise out-of-scope.

The LFS sample design is set up to yield about 60,000 households. Surveys needing smaller sample sizes usually "reserve" from one to six rotations per province, a rotation being one-sixth of the total sample. Sample stabilization is used to maintain the sample at desired levels, as when two rotations are reserved but the sample size needed only represents 1.5 rotations.

Requirements specific to the NPHS led to two modifications to this sampling strategy. The number of "reserves" needed was specified at the stratum level rather than the provincial level in order to meet the specific sub-provincial sample size requirements. It was also required that the number of clusters selected per stratum

be a multiple of four for variance estimation and seasonal representativity (this allowed strata to have two or more independent samples of four clusters each—one per collection period). As NPHS usually requested only between two and six clusters per LFS stratum, similar LFS strata were grouped to form larger NPHS strata with the required number of sample clusters.

As a result of these modifications, the NPHS sample of clusters can be considered as a stratified replicated sample where strata are groups of LFS strata and replicates are typically independent, identically distributed samples of four clusters each. There were exceptions, but they are not expected to have a significant impact on survey results.

### **5.2.5 Sample Design in Québec**

In Québec the NPHS sample is selected from dwellings participating in a health survey organized by Santé Québec: the 1992-93 *Enquête sociale et de santé* (ESS). The survey sampled 16,010 dwellings using a two-stage design similar to that of the LFS. The province was divided geographically by crossing 15 health areas with four urban density classes (Montréal Census Metropolitan Area, regional capitals, small urban agglomerations and the rural sector). In each area clusters were stratified by socio-economic characteristics and selected using a PPS sample. Selected clusters were enumerated and random samples of their dwellings were drawn: 10 per cluster in major cities, 20 or 30 elsewhere.

Santé Québec provided non-confidential information which allowed the classification of their sample into four types of households: one-member households; households with children; other households with youths (persons aged under 25); and the rest (more than one member and no youth or child). A household type was determined by NPHS personnel for the ESS non-respondents.

The NPHS sample size was first allocated among the four urban density classes. To avoid having too much sample in Montréal the allocation was proportional to  $\%(2W_h^2 + 1/4^2)$ , where  $W_h$  is the population share for class  $h$ ,  $h=1,2,3,4$ . In each class an attempt was made to obtain a sub-sample from the ESS which, as far as the selected panel member was concerned, would be proportional to the populations for the four household types. This was done by drawing a sufficient number of households from the ESS to give the required yield for households with children (the most underrepresented group), and then removing excess sample from the other three household groups. An initial sample which was almost 50% higher than needed was thus selected. After removing from it 2/3 of the one-member households, 1/2 of the other households with no youths or children, and 1/6 of households with youths but no children, the objective was nearly attained.

Considerations for seasonal representation and variance estimation, and integration with the NLSCY, affected the sub-sampling in Québec as they did elsewhere. ESS strata were thus collapsed to allow the formation of replicates, with the clusters in each replicate covering all four quarters (two quarters are covered per cluster in the rural and small urban sectors as sample sizes are higher there). The sample of households with children was split into an "Adult" sample and a "Children" sample by a 3:2 ratio, the terms having the same meaning as in other provinces. "Children" sample households in quarters 1 and 2 were reassigned to quarters 3 and 4. Since NPHS surveyed the current occupants of dwellings selected for the ESS, and changes occurred in some of those dwellings, the samples of households without children for quarters 3 and 4 were also to be split, by a 2:3 ratio, into an "Adult" and a "Children" sample.

**References:**

- <sup>1</sup> Kish, L. (1988). Multipurpose Sample Design, *Survey Methodology*, 14, 19-32.

## **6. Data Collection**

### **6.1 Questionnaire Design and Data Collection Method**

The NPHS questions were designed for computer-assisted interviewing (CAI), which meant that, as the questions were developed, the associated logical flow into and out of the questions was programmed. This included specifying the type of answer required, the minimum and maximum values, on-line edits associated with the question and what to do in case of item non-response.

With CAI, the interview can be controlled based on answers provided by the respondent. On-screen prompts are shown when an invalid entry is recorded and thus immediate feedback is given to the respondent and/or the interviewer to correct inconsistencies. Another enhancement is the automatic insertion of reference periods based on current dates. Prefilling of text or data based on information gathered during the interview allows the interviewer to proceed without having to search back for previous answers. This type of prefill includes such things as using the correct name or sex within the questions themselves. Allowable ranges/answers based on data collected during the interview can also be programmed. In other words the questionnaire can be customized to the respondent based on data collected at that time.

### **6.2 Tests**

A number of tests were conducted before the main survey was implemented in the field.

Focus groups were held during development stages of the questionnaires to study various aspects of their content. The main objectives were to determine the clarity and the quality of the questions, respondent reactions to sections that were felt to be sensitive (HIV, etc.), and to obtain approximate times for the length of the different sections.

Two field tests were also conducted. The tests involved four of Statistics Canada's Regional Offices. For the core sample, interviews were carried out by experienced Labour Force Survey interviewers. For the RDD sample, due to the volume increase, LFS interviewers were not used in all cases. The main objectives of the two tests were again to observe respondent reaction to the survey, to obtain estimates of time for the various sections, and to study the response rates. Field operations and procedures, interviewer training, and the computer program application (i.e., the questionnaire on computer) were also tested.

In addition to the field tests, the computer program application was extensively tested in-house in order to identify any errors in the program flow and text. Computer application testing was an ongoing operation up until the start of the main survey.



### **6.3 Interviewing**

Collection for the core sample was divided into four quarters (June, August and November 1996, and February 1997). The interviewers are part-time employees hired and trained specifically to carry out surveys using the computer-assisted interviewing method. An additional collection was held in June 1997 with further tracing attempts of non-respondents from previous quarters.

Collection for the RDD samples was carried out monthly, with survey start and end dates depending upon provincial funding. For Alberta, collection ran from June 1996 to March 1997. For Ontario, collection ran from October 1996 to August 1997. In Manitoba, collection ran from November 1996 to August 1997. Interviews were conducted by part-time employees, usually experienced in telephone interviewing, who were hired and trained specifically to carry out the RDD portion using CAI.

Respondents in the core sample were first contacted by telephone, and 95% of the interviews were done over the telephone. Personal visits were made if the respondent did not have a telephone, if the interviewer made a personal visit in the course of tracing a respondent, or upon request by the respondent. For the RDD collection, no personal interviews were allowed. The total interview took an average of one hour in each household.

In all dwellings, information about all household members was obtained from the person at home at the time of the interviewer call. Such "proxy" reporting, which accounted for approximately 55% of the information collected for this part of the interview, is used to reduce the cost of collection for the general component.

Proxy reporting of the health component was allowed for the selected respondent only for reasons of illness or incapacity. Such proxy reporting accounted for 2% of the information collected for respondents aged 12 years and older. On the other hand, all interviews for selected respondents under 12 years old were done by proxy.

### **6.4 Supervision and Control**

For the core sample, the LFS supervisory and control structure was employed for the NPHS collection. A similar structure of interviewers, senior interviewers and program managers was created for the RDD collection. Each RDD project was done in a central office, so managers were on-site to provide constant supervision.

## **6.5 Non-response to the NPHS**

Interviewers were instructed to make all reasonable attempts to obtain NPHS interviews with members of eligible households. For individuals in the core sample who at first refused to participate in the NPHS, a letter was sent from the Regional Office to the longitudinal respondent, stressing the importance of the survey and the household's cooperation. This was followed by a second call (or visit) from the interviewer. For cases in which the timing of the interviewer's call (or visit) was inconvenient, an appointment was made to call back at a more convenient time. If no one was home, numerous callbacks were made. For the RDD sample, several calls were made to non-response telephone numbers. When requested, a letter with information about the survey was sent from the Regional Office to the respondent, stressing the importance of the survey and the household's cooperation. Under no circumstances were sampled units replaced by other units for reasons of non-response.

## **6.6 Non-response Follow-up**

Many strategies were put in place to reduce the number of non-response cases. Before interviews started, a maximum recommended assignment size by interviewer was calculated based on test results. This allowed for the efficient follow-up of non-contact cases (i.e., to avoid overburdening interviewers).

Interviewer training covered ways of reducing the number of non-contacts (e.g., making calls or visits at various times of the day) using contact information given in the previous interview (for the core sample) and using telephone directory assistance to validate the number (for the RDD sample).

Refusals were followed up by senior interviewers, project supervisors or by other interviewers to try to convince respondents to participate in the survey.

To maximize the response rate, a large number of non-response cases were also followed up in subsequent collection periods.

## **6.7 Tracing**

For the core sample, the failure to trace a longitudinal respondent was an additional type of non-response. Interviewers had several ways to trace a respondent. The last known address and telephone number were provided as part of the information on the case, as well as the name and address of one or two previous contacts, if collected. In addition, interviewers were trained to follow up available leads such as local telephone directories and directory assistance. If these leads were unsuccessful, the case was transmitted to an experienced interviewer specially trained in tracing respondents. Tracer interviewers had access to Canada-wide telephone directories and reverse directories. The non-response rate due to failure to trace the longitudinal respondent was 1.7%, which is exceedingly low.

## **7. Data Processing**

### **7.1 Editing**

Editing was performed on-line in the computer-assisted interviewing (CAI) application during data collection. It was not possible to enter out-of-range values and flow errors were controlled through the use of CAI. Invalid values could not be entered and the correct question paths were automatically followed. For example, CAI ensured that questions that did not apply to the respondent were not asked. In other situations, warning messages were invoked, but no corrective action was taken (e.g., if an interviewer entered contradictory responses between questions). Because no corrective action was taken in such instances, edits were developed to be performed after data collection at Head Office. Inconsistencies were usually corrected by setting one or both of the variables in question to "not stated". No imputation was performed.

### **7.2 Coding**

Several questions allowing write-in responses had the write-in information coded into either new unique categories, or to a listed category if the write-in information duplicated a listed category. Where possible (e.g., occupation, industry, diseases), the coding followed the standard classification systems as used either in the Census of Population or in other Statistics Canada surveys such as the Health and Activity Limitation Survey and General Social Survey-cycle 6.

### **7.3 Creation of Derived and Grouped Variables**

To facilitate data analysis, a number of variables on the file have been derived using items found on the NPHS questionnaires. Derived variable names generally have a "D" or "G" in the fifth character of the variable name. In some cases, the derived variables are straightforward, involving collapsing of response categories. In other cases, several variables have been combined to create a new variable. Appendix F provides the details on how these more complex variables were derived.

### **7.4 Weighting**

The principle behind estimation in a probability sample such as the NPHS is that each person in the sample "represents", besides himself or herself, several other persons not in the sample. For example, in a simple random 2% sample of the population, each person in the sample represents 50 persons in the population. In the terminology used here, it can be said that each person has a weight of 50.

The weighting phase is a step which calculates, for each person, his or her associated weight. This weight appears on the microdata file, and must be used to derive meaningful estimates from the survey. For example, if the number of individuals who smoke daily (see question SMC6\_2 in section 9.2) is to be estimated, it is done by selecting the records referring to those individuals in the sample having that characteristic and summing the weights entered on those records.

Details of the method used to calculate these weights are presented in Chapter 11.

## **7.5 Suppression of Confidential Information**

It should be noted that the "public use" microdata files differ in a number of important respects from the survey "master" files held by Statistics Canada. These differences are the result of actions taken to protect the anonymity of individual survey respondents through suppression of individual values, variable grouping, and variable capping. Users requiring access to information excluded from the microdata files have two options—to purchase custom tabulations or to use the remote access option. Remote access allows computer programs to be submitted by users for processing at Statistics Canada. For more information on remote access see Chapter 12. Estimates generated will be released to the user, subject to meeting the guidelines for analysis and release outlined in Chapter 9 of this document.

**8. Data Quality**

**8.1 Response Rates**

The calculation of response rates for the NPHS must take into account the augmentation of the sample by RDD buy-in samples in Alberta, Manitoba, and Ontario, used for cross-sectional estimates only. Cross-sectional response rates are thus calculated separately for the core and for the RDD portion, and overall. The following table contains a summary of the 1996-97 response rates:

**1996-97 Response Rates Table**

Level	Core	RDD	Overall (Core+RDD)
Household	94.3%	80.0%	82.6%
Selected members : excluding RDD- Child	98.7%	94.8%	95.6%
Selected members : RDD-Child only		98.2%	98.2%

The following is a description of how the household response rate and the selected person response rates were calculated, for the cross-sectional and for the longitudinal files. It should be noted that out-of-scope dwellings (i.e., households not eligible for the sample) were not used in any of the calculations. Note that in the three sections below, the selected person response rate is calculated based on the number of responding households. To get an idea of the overall rate of response for selected person (based on all households), the two rates can be multiplied.

### **8.1.1 Core Cross-sectional Response Rates**

#### ***Household (HH) response rate***

$$\frac{\# \text{ of responding households}}{\text{all in-scope households}}$$

The 1996-97 *core cross-sectional* response rate is based on all continuing households, excluding those out of scope (e.g., a selected person who has moved out of the country). A responding household had *at least* one general-component questionnaire (H05) completed for a member of the household. This response rate at the Canada level for the NPHS was **94.3%**. At the provincial level, this rate varied from 91.8% in British Columbia to 96.5% in Newfoundland.

#### ***Selected person (SP) response rate***

The core selected person response rate can be thought of as the number of health-component questionnaires (H06) that *were* completed compared to the number that *should have been* completed.

$$\frac{\# \text{ of responding H06s}}{\# \text{ of responding households (i.e., eligible to answer)}}$$

For the core, the selected person response rate for the NPHS was **98.7%** at the Canada level, and ranged from 98.3% in Québec to 99.3% in New Brunswick and Alberta.

#### ***Relevant information for calculation of response rates:***

Number of respondents at the household level:	16,215
Number of respondents at the selected person level:	15,681
Number of non-respondents at the household level:	974
Number of non-respondents at the selected person level:	206
Number of out-of-scope households:	87

Calculation of household response rate:

$$\text{HH Rate} = \frac{16,215}{16,215 + 974} = \frac{16,215}{17,189} = 94.3\%$$

Calculation of selected person response rate:

$$\text{SP Rate} = \frac{15,681}{15,681 + 206} = \frac{15,681}{15,887} = 98.7\%$$

### **8.1.2 RDD Cross-sectional Response Rates (by province and total)**

#### ***Household response rate***

$$\frac{\# \text{ of responding households}}{\text{all in-scope households}}$$

For the RDD households only, the response rate for all three provinces combined was **80.0%**. At the provincial level, the rates were:

$$\begin{aligned} \text{Ontario} &= 37,795 / 48,770 = 77.5\% \\ \text{Manitoba} &= 11,210 / 12,952 = 86.6\% \\ \text{Alberta} &= 13,531 / 16,459 = 82.2\% \end{aligned}$$

For the RDD sample, out-of-scope households also included, for example, non-working or business telephone numbers.

#### ***Selected person response rate - Non-children (3 provinces)***

$$\frac{\# \text{ of responding non-child H06s}}{\# \text{ of responding households (i.e., eligible to answer)}}$$

The selected person response rate for non-children for all three provinces combined was **94.8%**. At the provincial level, the rates were:

$$\begin{aligned} \text{Ontario} &= 35,527 / 37,795 = 94.0\% \\ \text{Manitoba} &= 10,840 / 11,210 = 96.7\% \\ \text{Alberta} &= 12,925 / 15,531 = 95.5\% \end{aligned}$$



***Selected person response rate - Children (2 provinces only)***

$$\frac{\text{\# of responding child H06s}}{\text{\# of responding households that had a child <12 years old (i.e., eligible to answer)}}$$

The selected person response rate for children for Manitoba and Alberta combined was **98.2%**. At the provincial level, the rates were:

$$\text{Manitoba:} = 2,887 / 2,935 = 98.4\%$$

$$\text{Alberta:} = 3,944 / 4,020 = 98.1\%$$

***Relevant information for calculation of response rates:***

Number of respondents at the household level:	62,536
Number of non-child respondents at the selected person level:	59,292
Number of child respondents at the selected person level:	6,831
Number of non-respondents at the household level:	15,645
Number of non-child non-respondents at the selected person level:	3,244
Number of child non-respondents at the selected person level:	124
Number of out-of-scope households:	83,645

Calculation of household response rate:

$$\text{HH Rate} = \frac{62,536}{62,536 + 15,645} = \frac{62,536}{78,181} = 80.0\%$$

Calculation of non-child selected person response rate:

$$\text{SP Rate} = \frac{59,292}{59,292 + 3,244} = \frac{59,292}{62,536} = 94.8\%$$

Calculation of child selected person response rate:

$$\text{SP Rate} = \frac{6,831}{6,831 + 124} = \frac{6,831}{6,955} = 98.2\%$$

### **8.1.3 Overall Cross-sectional Response Rates**

#### ***Household response rate***

$$\frac{\# \text{ of responding households (Core + RDD)}}{\# \text{ of in-scope households (Core + RDD)}}$$

For the core and RDD households combined, the cross-sectional response rate at the Canada level was **82.6%**. At the provincial level, this rate varied from 78.8% in Ontario to 87.3% in Manitoba for those provinces with RDD samples and from 91.6% in British Columbia to 96.5% in Newfoundland for the others.

#### ***Selected person response rate - Excluding RDD children***

$$\frac{\# \text{ of responding non-RDD-child H06s (Core + RDD)}}{\# \text{ of responding households (Core + RDD)}}$$

The selected person response rate excluding RDD children was **95.6%** at the Canada level, and ranged from 94.4% in Ontario to 95.9% in Alberta for those with RDD samples and from 98.3% in Québec to 99.3% in New Brunswick for the others.

#### ***Selected person response rate - RDD children (2 provinces only)***

Note that this is the same calculation as in section 8.1.2 above.

$$\frac{\# \text{ of responding RDD child H06s}}{\# \text{ of responding households that had a child <12 years old (i.e., eligible to answer)}}$$

The selected person response rate for RDD children for Manitoba and Alberta combined was **98.2%**. At the provincial level, the rates were:

$$\begin{aligned} \text{Manitoba:} &= 2,887 / 2,935 = 98.4\% \\ \text{Alberta:} &= 3,944 / 4,020 = 98.1\% \end{aligned}$$

***1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION***

---

***Relevant information for calculation of response rates:***

Number of respondents at the household level:	78,751
Number of respondents at the selected person level, excluding RDD children:	74,973
Number of RDD child respondents at the selected person level:	6,831
Number of non-respondents at the household level:	16,619
Number of non-respondents at the selected person level, excluding RDD children:	3,450
Number of RDD child non-respondents at the selected person level:	124
Number of out-of-scope households:	83,732

Calculation of household response rate:

$$\text{HH Rate} = \frac{78,751}{78,751 + 16,619} = \frac{78,751}{95,370} = 82.6\%$$

Calculation of selected person response rate, excluding RDD children:

$$\text{SP Rate} = \frac{74,973}{74,973 + 3,450} = \frac{74,973}{78,423} = 95.6\%$$

Calculation of RDD child selected person response rate:

$$\text{SP Rate} = \frac{6,831}{6,831 + 124} = \frac{6,831}{6,955} = 98.2\%$$

## **8.2 Survey Errors**

The survey produces estimates based on information collected from a sample of individuals. Somewhat different estimates might have been obtained if a complete census had been taken using the same questionnaire, interviewers, supervisors, processing methods, etc. as those used in the survey. The difference between the estimates obtained from the sample and those resulting from a complete count taken under similar conditions is called the sampling error of the estimate.

Errors that are not related to sampling may occur at almost every phase of a survey operation. Interviewers may misunderstand instructions, respondents may make errors in answering questions, the answers may be incorrectly entered and errors may be introduced in the processing and tabulation of the data. These are all examples of non-sampling errors.

Over a large number of observations, randomly occurring errors will have little effect on estimates derived from the survey. However, errors occurring systematically will contribute to biases in the survey estimates. Considerable time and effort was made to reduce non-sampling errors in the survey. Quality assurance measures were implemented at each step of the data collection and processing cycle to monitor the quality of the data. These measures included the use of highly-skilled interviewers, extensive training of interviewers with respect to the survey procedures and questionnaire, observation of interviewers to detect problems and procedures to ensure that data collection errors were minimized.

A major source of non-sampling errors in surveys is the effect of non-response on the survey results. The extent of non-response varies from partial non-response (failure to answer just one or some questions) to total non-response. Partial non-response to NPHS was basically non-existent; once the questionnaire was started, it tended to be completed with very little non-response. Total non-response occurred because the interviewer was either unable to trace the respondent, no member of the household was able to provide the information, or the respondent refused to participate in the survey. Total non-response was handled by adjusting the weight of households that responded to the survey to compensate for those who did not respond.

In most cases, partial non-response to the survey occurred when the respondent did not understand or misinterpreted a question, refused to answer a question, could not recall the requested information, or could not provide personal or proxy information.

## **1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

---

This section of the documentation outlines the measures of sampling error that Statistics Canada commonly uses and that it urges users producing estimates from this microdata file to use also. Since it is an unavoidable fact that estimates from a sample survey are subject to sampling error, sound statistical practice calls for researchers to provide users with some indication of the magnitude of this sampling error.

The basis for measuring the potential size of sampling errors is the standard error of the estimates derived from survey results.

However, because of the large variety of estimates that can be produced from a survey, the standard error of an estimate is usually expressed relative to the estimate to which it pertains. This resulting measure, known as the coefficient of variation (C.V.) of an estimate, is obtained by dividing the standard error of the estimate by the estimate itself and is expressed as a percentage of the estimate.

For example, suppose that, based upon the survey results, one estimates that 24% of Canadians aged 12 and over are daily cigarettes smokers is found to have standard error of .003. Then the coefficient of variation of the estimate is calculated as:

$$\left( \frac{.003}{.24} \right) \times 100 \% = 1.25 \%$$

## **9. Guidelines For Tabulation, Analysis And Release**

This section of the documentation outlines the guidelines to be adhered to by users tabulating, analyzing, publishing or otherwise releasing any data derived from the survey microdata files. With the aid of these guidelines, users of microdata should be able to produce figures that are in close agreement with those produced by Statistics Canada and, at the same time, will be able to develop currently unpublished figures in a manner consistent with these established guidelines.

### **9.1 Rounding Guidelines**

In order that estimates for publication or other release derived from these microdata files correspond to those produced by Statistics Canada, users are urged to adhere to the following guidelines regarding the rounding of such estimates:

- a) Estimates in the main body of a statistical table are to be rounded to the nearest hundred units using the normal rounding technique. In normal rounding, if the first or only digit to be dropped is 0 to 4, the last digit to be retained is not changed. If the first or only digit to be dropped is 5 to 9, the last digit to be retained is raised by one. For example, in normal rounding to the nearest 100, if the last two digits are between 00 and 49, they are changed to 00 and the preceding digit (the hundreds digit) is left unchanged. If the last digits are between 50 and 99 they are changed to 00 and the preceding digit is incremented by 1.
- b) Marginal sub-totals and totals in statistical tables are to be derived from their corresponding unrounded components and then are to be rounded themselves to the nearest 100 units using normal rounding.
- c) Averages, proportions, rates and percentages are to be computed from unrounded components (i.e., numerators and/or denominators) and then are to be rounded themselves to one decimal using normal rounding. In normal rounding to a single digit, if the final or only digit to be dropped is 0 to 4, the last digit to be retained is not changed. If the first or only digit to be dropped is 5 to 9, the last digit to be retained is increased by 1.
- d) Sums and differences of aggregates (or ratios) are to be derived from their corresponding unrounded components and then are to be rounded themselves to the nearest 100 units (or the nearest one decimal) using normal rounding.
- e) In instances where, due to technical or other limitations, a rounding technique other than normal rounding is used resulting in estimates to be published or otherwise released that differ from corresponding estimates published by Statistics Canada, users are urged to note the reason for such differences in the publication or release document(s).

- f) Under no circumstances are unrounded estimates to be published or otherwise released by users. Unrounded estimates imply greater precision than actually exists.

## **9.2 Sample Weighting Guidelines for Tabulation**

The sample design used for the NPHS was not self-weighting. That is to say, the sampling weights are not identical for all individuals in the sample. When producing simple estimates, including the production of ordinary statistical tables, users must apply the proper sampling weight.

If proper weights are not used, the estimates derived from the microdata files cannot be considered to be representative of the survey population, and will not correspond to those produced by Statistics Canada.

Users should also note that some software packages may not allow the generation of estimates that exactly match those available from Statistics Canada, because of their treatment of the weight field.

### **9.2.1 Definitions of Types of Estimates: Categorical vs. Quantitative**

Before discussing how the NPHS data can be tabulated and analyzed, it is useful to describe the two main types of point estimates of population characteristics that can be generated from the microdata file for the National Population Health Survey.

#### Categorical Estimates:

Categorical estimates are estimates of the number, or percentage of the surveyed population possessing certain characteristics or falling into some defined category. The number of individuals who smoke daily is an example of such an estimate. An estimate of the number of persons possessing a certain characteristic may also be referred to as an estimate of an aggregate.

#### Example of Categorical Question:

SMK6\_2: At the present do/does ... smoke cigarettes daily, occasionally or not at all?

- Daily
- Occasionally
- Not at all

Quantitative Estimates:

Quantitative estimates are estimates of totals or of means, medians and other measures of central tendency of quantities based upon some or all of the members of the surveyed population.

An example of a quantitative estimate is the average number of cigarettes smoked per day by individuals who smoke daily. The numerator is an estimate of the total number of cigarettes smoked per day by individuals who smoke daily, and its denominator is an estimate of the number of individuals who smoke daily.

Example of Quantitative Question:

SMK6\_4: How many cigarettes do/does you/he/she smoke each day now?

|\_|\_| Number of cigarettes

**9.2.2 Tabulation of Categorical Estimates**

Estimates of the number of people with a certain characteristic can be obtained from the microdata file by summing the final weights of all records possessing the characteristic(s) of interest. Proportions and ratios of the form  $\hat{x} / \hat{y}$  are obtained by:

- a) summing the final weights of records having the characteristic of interest for the numerator ( $\hat{x}$ ),
- b) summing the final weights of records having the characteristic of interest for the denominator ( $\hat{y}$ ), then
- c) dividing the numerator estimate by the denominator estimate.

**9.2.3 Tabulation of Quantitative Estimates**

Estimates of quantities can be obtained from the microdata file by:

- a) multiplying the value of the variable of interest by the final weight and summing this quantity over all records of interest to obtain the numerator ( $\hat{x}$ ),
- b) summing the final weights of records having the characteristic of interest for the denominator ( $\hat{y}$ ), then
- c) dividing the numerator estimate by the denominator estimate.



For example, to obtain an estimate of the average number of cigarettes smoked each day by individuals who smoke daily, multiply the value reported in SMK6\_4 by the weight, WT66, then sum this value over those records with a response of "daily" to SMK6\_2 to obtain the numerator ( $\sum x_i$ ). Sum the final weight of those records with a response of "daily" to SMK6\_2 to obtain the denominator ( $\sum w_i$ ). Divide ( $\sum x_i$ ) by ( $\sum w_i$ ) to obtain the average number of cigarettes smoked each day by daily smokers.

### **9.3 Guidelines for Statistical Analysis**

The National Population Health Survey is based upon a complex design, with stratification and multiple stages of selection, and unequal probabilities of selection of respondents. Using data from such complex surveys presents problems to analysts because the survey design and the selection probabilities affect the estimation and variance calculation procedures that should be used.

While many analysis procedures found in statistical packages allow weights to be used, the meaning or definition of the weight in these procedures differs from that which is appropriate in a sample survey framework, with the result that while in many cases the estimates produced by the packages are correct, the variances that are calculated are almost meaningless.

For many analysis techniques (for example linear regression, logistic regression, analysis of variance), a method exists that can make the application of standard packages more meaningful. If the weights on the records are rescaled so that the average weight is one (1), then the results produced by the standard packages will be more reasonable; they still will not take into account the stratification and clustering of the sample's design, but they will take into account the unequal probabilities of selection. The rescaling can be accomplished by using in the analysis a weight equal to the original weight divided by the average of the original weights for the sampled units (people) contributing to the estimator in question.

In order to provide a means of assessing the quality of tabulated estimates, Statistics Canada has produced a set of Approximate Sampling Variability Tables (commonly referred to as "C.V. Tables") for the NPHS. These tables can be used to obtain approximate coefficients of variation for categorical-type estimates and proportions. See Chapter 10 for more details.

As an alternative to the use of the C.V. tables, a series of "bootstrap" files and associated programs are being supplied to users so that they will be able to calculate more precise individual variances to assess the quality of tabulated estimates. Again, see Chapter 10 for more details.

**9.4 Release Guidelines**

Before releasing and/or publishing any estimate from these microdata files, users should first 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 coefficient of variation for this estimate. For weighted estimates based on sample sizes of 30 or more, users should determine the coefficient of variation of the rounded estimate and follow the guidelines below.

**Sampling Variability Guidelines**

Type of Estimate	C.V. (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 recommends not to release 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 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. These data and any consequent findings should not be published. If the user chooses to publish these data or findings, then this disclaimer must be published with the data.”</p>



**10. Approximate Sampling Variability Tables**

In order to supply coefficients of variation that would be applicable to a wide variety of categorical estimates produced from this microdata file and that could be readily accessed by the user, a set of Approximate Sampling Variability Tables has been produced. These "look-up" tables allow the user to obtain an approximate coefficient of variation based on the size of the estimate calculated from the survey data.

The coefficients of variation (C.V.) are derived using the variance formula for simple random sampling and incorporating a factor which reflects the multi-stage, clustered nature of the sample design. This factor, known as the design effect, was determined by first calculating design effects for a wide range of characteristics and then choosing from among these a conservative value to be used in the look-up tables, which would then apply to the entire set of characteristics.

The six tables below show the design effects, sample sizes and population counts used to produce the six sets of Approximate Sampling Variability Tables. The six sets correspond to: the provincial and Canada levels for both household members and selected members; various age groups at the Canada level for both household members and selected members; and for Ontario, Manitoba and Alberta health area level for both household and selected members.

**Input Data for Provincial and Canada Level Sampling Variability Tables  
for Household Members (All Ages)**

<b>PROVINCE</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
Newfoundland	1.35	3,017	561,586
Prince Edward Island	1.53	2,752	135,717
Nova Scotia	1.41	2,775	917,152
New Brunswick	1.34	2,888	745,464
Québec	1.84	7,838	7,221,079
Ontario	1.45	110,845	11,131,489
Manitoba	3.67	32,399	1,085,634
Saskatchewan	1.35	2,785	974,046
Alberta	1.37	40,802	2,728,382
British Columbia	1.32	4,276	3,780,907
<b>CANADA</b>	5.26	210,377	29,281,456

**Input Data for Canada Level Age Group Sampling Variability Tables  
for Household Members (All Ages)**

<b>AGE GROUP</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
0-11	3.51	37,161	4,686,574
12-24	4.12	39,791	5,134,153
25-44	4.66	66,401	9,709,153
45-64	5.32	44,762	6,335,467
65+	6.01	22,262	3,416,109

**Input Data for Ontario, Manitoba and Alberta Health Area Level  
Sampling Variability Tables for Household Members (All Ages)**

<b>HEALTH AREA</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
<b>Ontario</b>			
Ottawa Carleton	1.13	7,143	767,447
Prescott,Russell,Stormont,Dundas,Glengarry,Renf.	1.19	3,222	287,829
Lanark,Leeds,Gren.,Hast.,P.E.,Fron.,Len.,Add.	1.20	5,247	497,455
Northumberland,Victoria,Haliburton,Peterborough	1.07	4,019	307,511
Durham	1.13	4,914	475,464
Peel	1.09	5,747	893,208
Metro Toronto	1.17	10,978	2,403,746
York	1.13	5,318	626,252
Simcoe	1.12	4,708	348,076
Halton	1.02	4,407	353,830
Niagara	1.02	4,041	418,650
Hamilton-Wentworth	1.09	4,606	484,533

**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

<b>HEALTH AREA</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
Wellington, Dufferin	1.13	4,504	226,271
Waterloo	1.06	4,393	423,839
Essex	1.08	4,437	354,704
Lambton, Kent	1.02	3,976	242,870
Elgin, Middlesex, Oxford	1.06	4,589	599,688
Bruce, Grey, Perth, Huron	1.03	4,298	301,559
Algoma, Cochrane	1.07	4,124	224,534
Manitoulin, Sudbury	1.00	3,953	205,646
Timiskaming, Muskoka, Parry Sound, Nipissing	1.06	3,852	222,474
Thunder Bay, Kenora, Rainy River	1.05	4,013	239,066
<b>Manitoba</b>			
North and South Eastman	1.27	6,676	86,532
Burntwood, Norman, Parkland	1.33	6,518	78,881
Central, Interlake	1.33	6,368	153,599
South Westman, Brandon, Marquette	1.20	8,501	111,448
Winnipeg	1.12	4,336	655,174
<b>Alberta</b>			
Northern Alberta	1.23	10,276	418,276
Southern Alberta	1.14	6,729	335,723
Central Alberta	1.14	8,719	427,650
Calgary	1.12	8,236	842,824
Edmonton	1.14	6,842	703,908

**Input Data for Provincial and Canada Level Sampling Variability Tables  
for Selected Members**

<b>PROVINCE</b>	<b>AGES</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
Newfoundland	2+	0.93	963	549,322
Prince Edward Island	2+	0.93	918	132,322
Nova Scotia	2+	0.94	986	895,914
New Brunswick	2+	0.92	1,032	728,118
Québec	2+	1.11	2,788	7,047,528
Ontario	2+	1.15	39,394	10,839,724
Manitoba	All	3.71	14,828	1,085,635
Saskatchewan	2+	0.93	1,047	948,511
Alberta (non-HPS questions)	All	1.60	18,305	2,728,383
British Columbia	2+	0.97	1,543	3,686,279
<b>CANADA</b>		4.32	81,804	28,641,735

Note: In Alberta, the RDD selected members were not asked certain HPS questions. The variables are listed in section 12.1.3. Therefore, when doing coefficient of variation calculations for these variables, use ages = 12+; design effect = 1.16; sample size = 1,278 and population = 2,243,982. A separate C.V. table is included in Appendix H for these variables.

**Input Data for Canada Level Age Group Sampling Variability Tables  
for Selected Members**

<b>AGE GROUP</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
2-11 (all provinces) <sup>1</sup>	8.18	7,224	3,925,489
0-11 (Man, Alta)	2.41	7,114	667,907
2-11 (other provinces)	1.93	1,288	3,378,942
12-24	5.59	12,120	5,134,153
25-44	4.14	28,900	9,709,158
45-64	3.88	19,019	6,335,468
65+	2.54	13,363	3,416,108

<sup>1</sup> Due to the sample design, the 2-11 age group for all of Canada has a rather large design effect.

**Input Data for Ontario, Manitoba and Alberta Health Area Level  
Sampling Variability Tables for Selected Members**

<b>HEALTH AREA</b>	<b>AGE</b>	<b>DESIGN EFFECT</b>	<b>SAMPLE SIZE</b>	<b>POPULATION</b>
<b>Ontario</b>	2+			
Ottawa Carleton		0.86	2,650	738,276
Prescott,Russell,Stormont,Dundas,Glengarry, Renf.		1.00	1,174	272,703
Lanark,Leeds,Grenville,Hast.,P.E.,Fron.,Len.,Add.		0.90	1,978	486,629
Northumberland,Victoria,Haliburton,Peterborough		1.00	1,513	277,441
Durham		0.72	1,618	482,391
Peel		0.83	1,728	882,317
Metro Toronto		1.05	4,085	2,238,057
York		0.84	1,582	607,098
Simcoe		0.82	1,642	339,940



**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

HEALTH AREA	AGE	DESIGN EFFECT	SAMPLE SIZE	POPULATION
Halton		0.77	1,493	360,691
Niagara		0.85	1,506	406,163
Hamilton-Wentworth		0.76	1,631	484,221
Brant, Haldiman, Norfolk		0.95	1,534	205,422
Wellington, Dufferin		0.63	1,546	250,593
Waterloo		0.90	1,555	383,946
Essex		0.87	1,558	339,105
Lambton, Kent		0.73	1,474	251,089
Elgin, Middlesex, Oxford		0.76	1,688	609,183
Bruce, Grey, Perth, Huron		0.86	1,573	297,485
Algoma, Cochrane		0.68	1,473	237,354
Manitoulin, Sudbury		0.77	1,479	206,157
Timiskaming, Muskoka, Parry Sound, Nipissing		0.70	1,465	233,708
Thunder Bay, Kenora, Rainy River		0.70	1,449	249,755
<b>Manitoba</b>	All			
North and South Eastman		1.29	2,953	86,530
Burntwood, Norman, Parkland		1.28	3,059	78,881
Central, Interlake		1.27	2,803	153,601
South Westman, Brandon, Marquette		1.17	4,050	111,449
Winnipeg		1.09	1,963	655,174
<b>Alberta</b>	All			
Northern Alberta		1.36	4,488	418,276
Southern Alberta		1.12	3,070	335,724
Central Alberta		1.09	3,927	427,650

**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

---

HEALTH AREA	AGE	DESIGN EFFECT	SAMPLE SIZE	POPULATION
Calgary		1.16	3,694	842,824
Edmonton		1.14	3,126	703,909

All coefficients of variation in the Approximate Sampling Variability Tables are *approximate* and, therefore, unofficial. Estimates of actual variance for specific variables may be obtained from Statistics Canada on a cost-recovery basis. As well, it is planned to supply "bootstrap" files and associated programs with this release. The use of actual variance estimates would allow users to release otherwise unreleaseable estimates, i.e., estimates with coefficients of variation in the "unacceptable" range.

Remember: 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 coefficient of variation for this estimate. This is because the formulas used for estimating the variance do not hold true for small sample sizes.

### **10.1 How to Use the C.V. Tables for Categorical Estimates**

The following rules should enable the user to determine the approximate coefficients of variation from the Sampling Variability Tables for estimates of the number, proportion or percentage of the surveyed population possessing a certain characteristic and for ratios and differences between such estimates.

#### **Rule 1: Estimates of Numbers Possessing a Characteristic (Aggregates)**

The coefficient of variation depends only on the size of the estimate itself. On the appropriate Sampling Variability Table, locate the estimated number in the left-most column of the table (headed "Numerator of Percentage") and follow the asterisks (if any) across to the first figure encountered. This figure is the approximate coefficient of variation.

#### **Rule 2: Estimates of Proportions or Percentages Possessing a Characteristic**

The coefficient of variation of an estimated proportion or percentage depends on both the size of the proportion or percentage and the size of the total upon which the proportion or percentage is based. Estimated proportions or percentages are relatively more reliable than the corresponding estimates of the numerator of the proportion or percentage, when the proportion or percentage is based upon a sub-group of the population. This is due to the fact that the

coefficients of variation of the latter type of estimates are based on the largest entry in a row of a particular table, whereas the coefficients of variation of the former type of estimators are based on some entry (not necessarily the largest) in that same row. (Note that in the tables the C.V.'s decline in value reading across a row from left to right).

For example, the estimated proportion of individuals who smoke daily out of those who smoke at all is more reliable than the estimated number who smoke daily.

When the proportion or percentage is based upon the total population covered by each specific table, the C.V. of the proportion or percentage is the same as the C.V. of the numerator of the proportion or percentage. In this case, Rule 1 can be used.

When the proportion or percentage is based upon a subset of the total population (e.g., those who smoke at all), reference should be made to the proportion or percentage (across the top of the table) and to the numerator of the proportion or percentage (down the left side of the table). The intersection of the appropriate row and column gives the coefficient of variation.

**Rule 3: Estimates of Differences Between Aggregates or Percentages**

The standard error of a difference between two estimates is approximately equal to the square root of the sum of squares of each standard error considered separately. That is, the standard error of a difference ( $\hat{d} = \hat{x}_2 - \hat{x}_1$ ) is:

$$s_{\hat{d}} = \sqrt{(\hat{x}_1 a_1)^2 + (\hat{x}_2 a_2)^2}$$

where  $\hat{x}_1$  is estimate 1,  $\hat{x}_2$  is estimate 2, and  $a_1$  and  $a_2$  are the coefficients of variation of  $\hat{x}_1$  and  $\hat{x}_2$  respectively. The coefficient of variation of  $\hat{d}$  is given by  $s_{\hat{d}} / \hat{d}$ . This formula is accurate for the difference between separate and uncorrelated characteristics, but is only approximate otherwise.

**Rule 4: Estimates of Ratios**

In the case where the numerator is a subset of the denominator, the ratio should be converted to a percentage and Rule 2 applied. This would apply, for example, to the case where the denominator is the number of individuals who smoke at all and the numerator is the number of individuals who smoke daily out of those who smoke at all.

Consider the case where the numerator is not a subset of the denominator, as for example, the ratio of the number of individuals who smoke daily or occasionally as compared to the number of individuals who do not smoke at all. The standard deviation of the ratio of the estimates is approximately equal to the square root of the sum of squares of each coefficient of variation considered separately multiplied by  $\hat{R}$ , where  $\hat{R}$  is the ratio of the estimates ( $\hat{R} = \hat{x}_1 / \hat{x}_2$ ). That is, the standard error of a ratio is:

$$s_{\hat{R}} = \hat{R} \sqrt{a_1^2 + a_2^2}$$

where  $a_1$  and  $a_2$  are the coefficients of variation of  $\hat{x}_1$  and  $\hat{x}_2$  respectively.

The coefficient of variation of  $\hat{R}$  is given by  $s_{\hat{R}} / \hat{R} = \sqrt{a_1^2 + a_2^2}$ . The formula will tend to overstate the error, if  $\hat{x}_1$  and  $\hat{x}_2$  are positively correlated and understate the error if  $\hat{x}_1$  and  $\hat{x}_2$  are negatively correlated.

**Rule 5: Estimates of Differences of Ratios**

In this case, Rules 3 and 4 are combined. The C.V.'s for the two ratios are first determined using Rule 4, and then the C.V. of their difference is found using Rule 3.

**10.2 Examples of Using the C.V. Tables for Categorical Estimates**

The following "real life" examples are included to assist users in applying the foregoing rules.

**Example 1 : Estimates of Numbers Possessing a Characteristic (Aggregates)**

Suppose that a user estimates that 5,794,518 individuals smoke daily in Canada. How does the user determine the coefficient of variation of this estimate?

- 1) Refer to the CANADA level C.V. table for SELECTED MEMBERS.
- 2) The estimated aggregate (5,794,518) does not appear in the left-hand column (the "Numerator of Percentage" column), so it is necessary to use the figure closest to it, namely 6,000,000.
- 3) The coefficient of variation for an estimated aggregate (expressed as a percentage) is found by referring to the first non-asterisk entry on that row, namely, 1.4%.
- 4) So the approximate coefficient of variation of the estimate is 1.4%. The finding that there were 5,794,518 individuals who smoke daily is publishable with no qualifications.

**Example 2 : Estimates of Proportions or Percentages Possessing a Characteristic**

Suppose that the user estimates that  $5,794,518/6,781,835=85.4\%$  of individuals in Canada who smoke at all smoke daily. How does the user determine the coefficient of variation of this estimate?

- 1) Refer to the CANADA level C.V. table for SELECTED MEMBERS.
- 2) Because the estimate is a percentage which is based on a subset of the total population (i.e., individuals who smoke at all, that is to say, daily or occasionally), it is necessary to use both the percentage (85.4%) and the numerator portion of the percentage (5,794,518) in determining the coefficient of variation.
- 3) The numerator, 5,794,518, does not appear in the left-hand column (the "Numerator of Percentage" column) so it is necessary to use the figure closest to it, namely 6,000,000. Similarly, the percentage estimate does not appear as any of the column headings, so it is necessary to use the figure closest to it, 90.0%.
- 4) The figure at the intersection of the row and column used, namely 0.5% is the coefficient of variation (expressed as a percentage) to be used.

- 5) So the approximate coefficient of variation of the estimate is 0.5%. The finding that 85.4% of individuals who smoke at all smoke daily can be published with no qualifications.

**Example 3 : Estimates of Differences Between Aggregates or Percentages**

Suppose that a user estimates that  $4,902,055/5,794,518=85\%$  of those who smoke daily smoke 10 or more cigarettes daily (estimate 1) while  $4,261,353/5,692,300=75\%$  of those who smoke occasionally or not at all, but at one time smoked daily, smoked 10 or more cigarettes daily at that time (estimate 2). Note that these estimates are based on the results of questions SMC6\_2, SMC6\_4, SMC6\_4A, SMC6\_5 and SMC6\_7. How does the user determine the coefficient of variation of the difference between these two estimates?

- 1) Using the CANADA level C.V. table for SELECTED MEMBERS in the same manner as described in example 2 gives the C.V. for estimate 1 as 0.5% (expressed as a percentage), and the C.V. for estimate 2 as 1.1% (expressed as a percentage).
- 2) Using rule 3, the standard error of a difference ( $\hat{d} = \hat{X}_2 - \hat{X}_1$ ) is:

$$s_{\hat{d}} = \sqrt{(\hat{X}_1 a_1)^2 \% (\hat{X}_2 a_2)^2}$$

where  $\hat{X}_1$  is estimate 1,  $\hat{X}_2$  is estimate 2, and  $a_1$  and  $a_2$  are the coefficients of variation of  $\hat{X}_1$  and  $\hat{X}_2$  respectively.

That is, the standard error of the difference  $\hat{d} = (.75 - .85) = .10$  is:

$$s_{\hat{d}} = \sqrt{[(.85)(.005)]^2 \% [(.75)(.011)]^2} = .009$$

- 3) The coefficient of variation of  $\hat{d}$  is given by  $s_{\hat{d}} / \hat{d} = .009/.10 = 0.095$ .
- 4) So the approximate coefficient of variation of the difference between the estimates is 9.5% (expressed as a percentage). This estimate can be published with no qualifications.

**Example 4 : Estimates of Ratios**

Suppose that the user estimates that 5,794,518 individuals smoke daily, while 987,317 individuals smoke occasionally. The user is interested in comparing the estimate of daily to occasional smokers in the form of a ratio. How does the user determine the coefficient of variation of this estimate?

- 1) First of all, this estimate is a ratio estimate, where the numerator of the estimate ( $= \hat{X}_1$ ) is the number of individuals who smoke occasionally. The denominator of the estimate ( $= \hat{X}_2$ ) is the number of individuals who smoke daily.
- 2) Refer to the CANADA level C.V. table for SELECTED MEMBERS.
- 3) The numerator of this ratio estimate is 987,317. The figure closest to it is 1,000,000. The coefficient of variation for this estimate (expressed as a percentage) is found by referring to the first non-asterisk entry on that row, namely, 3.8%.
- 4) The denominator of this ratio estimate is 5,794,518. The figure closest to it is 6,000,000. The coefficient of variation for this estimate (expressed as a percentage) is found by referring to the first non-asterisk entry on that row, namely, 1.4%.
- 5) So the approximate coefficient of variation of the ratio estimate is given by rule 4, which is,

$$a_{\hat{R}} = \sqrt{a_1^2 \% a_2^2}$$

where  $a_1$  and  $a_2$  are the coefficients of variation of  $\hat{X}_1$  and  $\hat{X}_2$  respectively.

That is ,

$$a_{\hat{R}} = \sqrt{(.038)^2 \% (.014)^2}$$

$$= 0.040$$

The obtained ratio of occasional to daily smokers is 987,317/5,794,518 which is 0.17:1. The coefficient of variation of this estimate is 4.0% (expressed as a percentage), which is releasable with no qualifications.

### **10.3 How to Use the C.V. Tables 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 the population is repeated indefinitely, 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 standard error of an estimate, confidence intervals for estimates may be obtained under the assumption that under repeated sampling of the population, the various estimates obtained for a population characteristic are normally distributed about the true population value. Under 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,  $\hat{X}$ , are generally expressed as two numbers, one below the estimate and one above the estimate, as  $(\hat{X}-k, \hat{X}+k)$  where  $k$  is determined depending upon the level of confidence desired and the sampling error of the estimate.

Confidence intervals for an estimate can be calculated directly from the Approximate Sampling Variability Tables by first determining from the appropriate table the coefficient of variation of the estimate  $\hat{X}$ , and then using the following formula to convert to a confidence interval CI:

$$CI_X = [\hat{X} - z \hat{X} a_{\hat{X}}, \hat{X} + z \hat{X} a_{\hat{X}}]$$

where  $a_{\hat{X}}$  is the determined coefficient of variation of  $\hat{X}$ , and

- $z = 1$  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.



#### 10.4 Example of Using the C.V. Tables to Obtain Confidence Limits

A 95% confidence interval for the estimated proportion of individuals who smoke daily from those who smoke at all (from example 2, section 10.2) would be calculated as follows.

$$\hat{X} = .854$$

$$z = 2$$

$a_x = .005$  is the coefficient of variation of this estimate as determined from the tables

$$CI_x = \{.854 - (2) (.854) (.005), .854 + (2) (.854) (.005)\}$$

$$CI_x = \{.846, .863\}$$

#### 10.5 How to Use the C.V. Tables to do a Z-test

Standard errors may also be used to perform hypothesis testing, a procedure for distinguishing between population parameters using sample estimates. The sample estimates can be numbers, averages, percentages, 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  $X_1$  and  $X_2$  be sample estimates for 2 characteristics of interest. Let the standard error on the difference  $\hat{X}_1 - \hat{X}_2$  be  $s_d$ .

If  $z = (\hat{x}_1 - \hat{x}_2) / s_d$  is between -2 and 2, then no conclusion about the difference between the characteristics is justified at the 5% level of significance. If however, this ratio is smaller than -2 or larger than +2, the observed difference is significant at the 0.05 level.

#### 10.6 Example of Using the C.V. Tables to do a Z-test

Let us suppose we wish to test, at 5% level of significance, the hypothesis that there is no difference between the proportion of individuals who smoke daily at a rate of 10 or more cigarettes AND the proportion of those who smoke occasionally or not at all, but at one time smoked daily at a rate of 10 or more cigarettes. From example 3, section 10.2, the standard error of the difference between these two estimates was found to be = .009.

Hence ,

$$z = \frac{\hat{X}_1 - \hat{X}_2}{s_{\hat{d}}} = \frac{.85 - .75}{.009} = \frac{.10}{.009} = 10.52$$

Since  $z = 10.52$  is greater than 2, it must be concluded that there is a significant difference between the two estimates at the 0.05 level of significance.

### **10.7 Exact Variances/Coefficients of Variation**

All coefficients of variation in the Approximate Sampling Variability Tables (C.V. Tables) are indeed approximate and, therefore, unofficial. However, exact coefficients of variation for specific variables may be obtained from Statistics Canada on a cost-recovery basis. As well, it is planned to supply "bootstrap" files and associated programs with this release. The types of estimates supported include aggregates, proportions, ratios, differences between aggregates, as well as more sophisticated types of analyses such as estimates of coefficients from linear regressions and logistic regressions, among others. The exact coefficients of variation are obtained via an exact variance program, which uses a technique called "bootstrapping". This technique involves dividing the records on the microdata files into subgroups (or replicates) and determining the variation in the estimates from replicate to replicate. There are a number of reasons why a user may require an exact variance. A few are given below.

Firstly, if a user desires estimates at a geographic level smaller than the province (for example, at the urban/rural level), then the C.V. tables provided are not adequate. Coefficients of variation of these estimates may be obtained using "domain" estimation techniques through the exact variance program.

Secondly, should a user require more sophisticated analyses such as estimates of coefficients from linear regressions or logistic regressions, the C.V. tables will not provide correct associated coefficients of variation. Although some standard statistical packages allow sampling weights to be incorporated in the analyses, the variances that are produced often do not take into account the stratified and clustered nature of the design properly, whereas the exact variance program would do so.

Thirdly, for estimates of quantitative variables, separate tables are required to determine their sampling error. Since most of the variables for the National Population Health Survey are primarily categorical in nature, this has not been done. Thus, users wishing to obtain coefficients of variation for quantitative variables can do so through the exact variance program. As a general rule, however, the coefficient of variation of a quantitative total will be larger than the coefficient of variation of the corresponding category estimate (i.e., the estimate of the number of persons contributing to the quantitative estimate). If the corresponding category estimate is not releasable, the

quantitative estimate will not be either. For example, the coefficient of variation of the estimate of the total number of cigarettes smoked each day by individuals who smoke daily would be greater than the coefficient of variation of the corresponding estimate of the number of individuals who smoke daily. Hence if the coefficient of variation of the latter is not releasable, then the coefficient of variation of the corresponding quantitative estimate will also not be releasable.

Lastly, should a user find himself/herself in a position where he/she can use the C.V. tables, but this renders a coefficient of variation in the "marginal" range (16.6% - 33.3%), the user should release the associated estimate with a warning cautioning users of the high sampling variability associated with the estimate. This would be a good opportunity to recalculate the coefficient of variation through the exact variance program to find out if it is releasable without a qualifying note. The reason for this is that the coefficients of variation produced by the tables are based on a wide range of variables and are therefore considered crude, whereas the exact variance program would give an exact coefficient of variation associated with the variable in question.

The exact variance/coefficient of variation program is available now and any user interested in this service should contact Stéphane Tremblay (613-951-4765) from the Health Statistics Methods Section within Household Survey Methods Division at Statistics Canada. Although there will be no charge for any computer time required, there will be a fee charged for any consultation time required to set up the request as well as for any time required to set up the associated computer runs. The daily consultation rate, based on a 7.5 hour day, is \$497.52; this rate may be broken down into an appropriate number of hours or minutes, if required. Naturally, the length of the consultation will vary from request to request and will depend upon the complexity of the analysis, the number of variables to be analyzed, etc. We also plan to produce a secondary release for the public use microdata file with "bootstrap" factors, so that users can calculate their own approximate coefficients of variation for most analyses.

### **10.8 Release Cut-offs for the NPHS**

The minimum cut-offs for estimates of totals at the provincial and Canada levels, those for various age groups at the Canada level and for Ontario, Manitoba and Alberta health areas, for both household members and selected members, are specified in the six tables below. Estimate sizes smaller than the minimum given in the "Unacceptable" column may not be released under any circumstances.

**Table of Release Cut-offs for Totals Based on Provincial/Canada Level Estimates for Household Members (All Ages)**

<b>PROVINCE</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
Newfoundland	9,000	4,000	2,500
Prince Edward Island	2,500	1,000	500
Nova Scotia	17,000	7,500	4,000
New Brunswick	12,500	5,500	3,000
Québec	61,500	27,000	15,500
Ontario	5,500	2,500	1,500
Manitoba	4,500	2,000	1,000
Saskatchewan	17,000	7,500	4,000
Alberta	3,500	1,500	1,000
British Columbia	42,500	18,500	10,500
<b>CANADA</b>	27,000	11,500	6,500

**Table of Release Cut-offs for Totals Based on Age Group Estimates at the Canada Level for Household Members (All Ages)**

<b>AGE GROUP</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
0-11	16,000	7,000	4,000
12-24	19,500	8,500	5,000
25-44	25,000	11,000	6,000
45-64	27,500	12,000	7,000
65+	33,500	14,500	8,500

**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

**Table of Release Cut-offs for Totals Based on Ontario, Manitoba and Alberta Health Area Level Estimates for Household Members (All Ages)**

HEALTH AREA	ACCEPTABLE	MARGINAL	UNACCEPTABLE
<b>Ontario</b>			
Ottawa Carleton	4,500	2,000	1,000
Prescott,Russell,Stormont,Dundas,Glengarry,Renf	4,000	1,500	1,000
Lanark,Leeds,Grenville,Hast.,P.E.,Fron.,Len.,Add.	4,000	2,000	1,000
Northumberland,Victoria, Haliburton, Peterb.	3,000	1,500	500
Durham	4,000	1,500	1,000
Peel	6,000	2,500	1,500
Metro Toronto	9,500	4,000	2,500
York	5,000	2,000	1,000
Simcoe	3,000	1,500	500
Halton	3,000	1,500	500
Niagara	4,000	1,500	1,000
Hamilton-Wentworth	4,000	2,000	1,000
Brant, Haldiman, Norfolk	2,000	1,000	500
Wellington, Dufferin	2,000	1,000	500
Waterloo	3,500	1,500	1,000
Essex	3,000	1,500	1,000
Lambton, Kent	2,500	1,000	500
Elgin, Middlesex, Oxford	5,000	2,000	1,000
Bruce, Grey, Perth, Huron	2,500	1,000	500
Algoma, Cochrane	2,000	1,000	500
Manitoulin, Sudbury	2,000	1,000	500
Timiskaming,Muskoka,Parry Sound,Nipissing	2,000	1,000	500

**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

<b>HEALTH AREA</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
Thunder Bay, Kenora, Rainy River	2,500	1,000	500
<b>Manitoba</b>			
North and South Eastman	500	500	500
Burntwood, Norman, Parkland	500	500	500
Central, Interlake	1,000	500	500
South Westman, Brandon, Marquette	500	500	500
Winnipeg	6,000	2,500	1,500
<b>Alberta</b>			
Northern Alberta	2,000	1,000	500
Southern Alberta	2,000	1,000	500
Central Alberta	2,000	1,000	500
Calgary	4,000	2,000	1,000
Edmonton	4,500	2,000	1,000

**Table of Release Cut-offs for Totals Based on Provincial/Canada  
Level Estimates for Selected Members**

<b>PROVINCE</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
Newfoundland	19,000	8,500	4,500
Prince Edward Island	4,500	2,000	1,000
Nova Scotia	30,500	13,500	7,500
New Brunswick	23,000	10,000	6,000
Québec	101,500	44,500	25,000
Ontario	11,500	5,000	3,000
Manitoba	10,000	4,500	2,500
Saskatchewan	30,000	13,500	7,500
Alberta	8,500	4,000	2,000
British Columbia	83,000	36,500	21,000
CANADA	55,500	24,000	13,500

**Table of Release Cut-offs for Totals Based on Age Group  
Estimates at the Canada Level for Selected Members**

<b>AGE GROUP</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
2-11 (all provinces)	156,500	70,000	39,500
0-11 (Man, Alta)	8,000	3,500	2,000
2-11 (other provinces)	176,500	79,000	45,000
12-24	85,500	37,500	21,500
25-44	51,000	22,000	12,500
45-64	47,000	20,500	11,500
65+	23,500	10,500	6,000

**Table of Release Cut-offs for Totals Based on Ontario, Manitoba and Alberta  
Health Area Level Estimates for Selected Members**

<b>HEALTH AREA</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
<b>Ontario</b>			
Ottawa Carleton	8,500	4,000	2,000
Prescott,Russell,Stormont,Dundas,Gleng.,Renf.	8,500	3,500	2,000
Lanark,Leeds,Gren.,Hast.,P.E.,Fron.,Len.,Add.	8,000	3,500	2,000
N'umberland, Vict., Halib., Peterb.	6,500	3,000	1,500
Durham	8,000	3,500	2,000
Peel	15,500	6,500	4,000
Metro Toronto	21,000	9,000	5,000
York	11,500	5,000	3,000
Simcoe	6,000	2,500	1,500
Halton	6,500	3,000	1,500



**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

<b>HEALTH AREA</b>	<b>ACCEPTABLE</b>	<b>MARGINAL</b>	<b>UNACCEPTABLE</b>
Niagara	8,000	3,500	2,000
Hamilton-Wentworth	8,000	3,500	2,000
Brant, Haldiman, Norfolk	4,500	2,000	1,000
Wellington, Dufferin	3,500	1,500	1,000
Waterloo	8,000	3,500	2,000
Essex	7,000	3,000	1,500
Lambton, Kent	4,500	2,000	1,000
Elgin, Middlesex, Oxford	10,000	4,500	2,500
Bruce, Grey, Perth, Huron	6,000	2,500	1,500
Algoma, Cochrane	4,000	1,500	1,000
Manitoulin, Sudbury	4,000	1,500	1,000
Timiskaming, Muskoka, Parry Sound, Nipissing	4,000	2,000	1,000
Thunder Bay, Kenora, Rainy River	4,500	2,000	1,000
<b>Manitoba</b>			
North and South Eastman	1,500	500	500
Burntwood, Norman, Parkland	1,000	500	500
Central, Interlake	2,500	1,000	500
South Westman, Brandon, Marquette	1,000	500	500
Winnipeg	13,000	6,000	3,500
<b>Alberta</b>			
Northern Alberta	4,500	2,000	1,000
Southern Alberta	4,500	2,000	1,000
Central Alberta	4,500	2,000	1,000
Calgary	9,500	4,000	2,500
Edmonton	9,500	4,000	2,500

## **11. Weighting**

The household component of the National Population Health Survey in 1994 had two basic designs, one for the nine provinces outside Québec, and one for Québec. In the nine provinces outside Québec, the NPHS used the design of the Labour Force Survey (LFS) with many modifications, to generate a sample of its own. Consequently, the derivation of weights was tied to the weighting procedure used for the LFS. In Québec, however, a two-phase sample design was implemented, where the first phase was drawn by the *Enquête sociale et de santé* (ESS) in 1992-93, and the second phase sample was drawn by the NPHS. Thus, in Québec, the derivation of the weights was tied to the weighting procedure used by the ESS. See Chapter 5 - Sample Design for more details. In 1996-97 additional independent samples were drawn in Ontario, Manitoba and Alberta using Random Digit Dialing (RDD) to allow for the production of reliable estimates at the health-area level.

In cycle 1, two sets of weights were required, one for cross-sectional estimates based upon variables from the general component of the questionnaire administered in 1994-95 to all household members, and one for cross-sectional estimates based on variables from the health component of the questionnaire administered to only the selected member. In cycle 2, three sets of weight are required for cross-sectional purposes. The first two are similar to those of 1994-95, but are calculated differently to take into account the sample design of 1996-97. The third weight is applicable only to the health file and is required for analysis of specific populations and variables. See Chapter 12 for more information.

The 1996-97 weighting procedure is based significantly on that of 1994-95. The 1994-95 final weights for each continuing selected respondent are taken as a starting point. Some weight adjustments that are no longer necessary in 1996-97 have been removed, to create a “stripped” weight for each selected respondent. From this point, the new adjustments are made to come up with the 1996-97 final weight, for both the selected (longitudinal) member and every member of the current household. In the following sections, a brief description of the 1994-95 weighting procedures still relevant for weighting in 1996-97 is included. A full description of the 1994-95 procedures may be found in the 1994-95 National Population Health Survey Public Use Microdata File documentation. Section 11.1 describes the new adjustments necessary in 1996-97 for the continuing (core) NPHS sample, while Section 11.2 describes the 1996-97 procedures in the provinces with RDD buy-ins (Ontario, Manitoba, and Alberta). Section 11.3 and Section 11.4 describe the 1994-95-based procedures for the provinces outside Québec and for Québec, respectively.

## **11.1 Cross-sectional Weighting for the 1996-97 NPHS—Core Household Sample**

This section describes the 1996-97 weighting procedures for selected members and all members of their households in the continuing (core) NPHS sample. A complete description of the additional weighting procedures necessary in the buy-in provinces (Ontario, Manitoba, and Alberta) is found in Section 11.2.

### **11.1.1 Stripped Weights**

As described in Sections 11.3 and 11.4, the starting point of the 1996-97 weighting procedure is the “stripped” weight, based upon the original sample design of 1994-95. Once these weights are obtained, the following weight adjustments are performed.

### **11.1.2 Weight Adjustments for Household Members**

#### Adjustment 1A: Household Non-response Adjustment

The definition of a non-responding household encompasses any of the following situations: refusal, special circumstance, language barrier, no one at home, temporarily absent or computer problem. There are also cases where it was determined that the selected member being followed-up in 1996-97 was dead, institutionalized, had moved to the Yukon or Northwest Territories, or out of the country. For cross-sectional weighting purposes, these households are included as responding households at this stage, but are subsequently dropped from further calculations. These units do not appear on the cross-sectional microdata file but do appear on the longitudinal file.

To adjust for cases of entire households that did not respond to the 1996-97 survey, the following adjustment is made:

$$\frac{\text{sum of weights for responding and non \& responding households within weighting class}}{\text{sum of weights for responding households within weighting class}}$$

Weighting classes consist of groupings of units (or households) that share the same propensity to respond to the survey. Characteristics from cycle 1, available for cycle 2 respondents and non-respondents alike, are used to define membership in the weighting classes. Classes are formed using a clustering algorithm that arranges the sample units into a tree structure by successively splitting the data set into “branches” based on the units’ characteristics. Each split aims to divide the units present into two or more groups that are most

dissimilar with respect to their observed non-response rate (and within which the non-response rates are expected to be more similar). A different characteristic may be used to define each split. For example, units may first be divided into owner-occupied dwellings and rented dwellings. The former split may then be further split into five groups based on the level of household income while the latter may be further split based on the respondent's age. Each of the newly formed groups may further be split, based on other characteristics, and so on. The results of the final splits are the weighting classes.

The software *Knowledge Seeker IV for Windows*, developed by ANGOSS Software International Limited, is used to generate the tree structure. We used an improved version of the CHAID (Chi-Square Automatic Interaction Detector) algorithm available in Knowledge Seeker to identify at each node the characteristic that best splits the sample into groups that are dissimilar with respect to a certain characteristic, here the response/non-response indicator.

When categorical data with more than two levels are used, Knowledge Seeker may group together one or more levels so that the number of "branches" may be less than the number of categories. For continuous characteristics, such as age, Knowledge Seeker first divides the data into ten ranges, which may or may not be collapsed, sometimes resulting in only two "branches". Statistical tests are performed at each step to ensure that only statistically-significant splits are generated. Bonferroni adjustments are made to the significance level of the individual tests to ensure that the significance level of each split is attained. The splitting ends when no more statistically-significant splits are found or when splits generate classes that are too small (a minimum of 30 units per class is used). For more information about the CHAID algorithm, see Kass (1980)<sup>3</sup>.

Personal characteristics of the selected respondent, as well as dwelling or household characteristics, are used to define the weighting classes for household non-response. The selected respondents' personal characteristics are deemed to play a significant role in predicting household non-response for several reasons. Often, person and household level non-response are equivalent. An obvious example is the fact that selected respondents who could not be located in 1996-97 led to a household non-response situation. Also, during data collection, emphasis was placed on getting a response for the selected respondent since information on the selected respondent was essential for longitudinal purposes. If the selected respondent was not available or did not respond, the interviewers were instructed not to complete the general component for the rest of the members of the household.

Finally, in many households, members shared common characteristics, such as race. In those cases, the respondent characteristics are in some sense also household characteristics.

Separate sets of weighting adjustment cells are created for each province. In addition to household and personal characteristics of the selected member, some characteristics that are related to the design of the survey were used, to reduce the effect of the sample design on the results of the statistical inference and incorporate design variables into the analysis. The characteristics vary by province but include the following variables:

Geographic	Province, Census Metropolitan Area, urban/rural indicator
Household	Dwelling type, owner/renter status, family type, household income adequacy, main source of income, non-response flag for income in 1994-95, presence of children in the household
Personal	Sex, age, age over 16 indicator, marital status, race, country of birth, age at immigration, restriction of activity flag, main activity/labour force status

Adjustment 2A: Interprovincial Migrations

It is sometimes necessary to make an adjustment for panel members that move from provinces with large populations to those with small populations, as the members' weights are often atypically large compared with those of other similar units in their new province. Corrections are only made in instances where at least one extreme weight was generated within the provincial move pattern in question. In those cases, the weights of all those selected members falling within that move pattern are reduced so that the sum of their weights equal the demographic projection of the number of movers within the move pattern in the past two years. If the adjustment would have created an increase, the adjustment was not made, since that would augment the extreme weight even further.

Adjustment 3A: Weight Share Method

Only selected (longitudinal) members are traced between cycles. Therefore if the composition of the household of the longitudinal member changes between cycles, their cohabitants at cycle 2 will not have any weights associated with them since they were not selected in the panel in cycle 1. The weight share method is a mechanism for assigning these individuals weights in such a way that resulting estimates are unbiased. Typically, every household member is assigned

the panel member's weight divided by the number of household members who were in-scope to the survey in cycle 1 (e.g., excluding persons born or entering Canada since 1994). For more information see Ernst (1989)<sup>2</sup> or Lavallée (1995)<sup>4</sup>.

**Adjustment 4A: Household Member Non-response**

Non-response at cycle 2 attributable to a household member is less than 2 percent. Therefore, all age/sex groups are combined and the following adjustment is made, by province:

$$\frac{\text{sum of weights of respondent and non \& respondent individuals in a province}}{\text{sum of weights of respondent individuals in a province}}$$

**Adjustment 5A: Integration of Core and Buy-in Samples**

As mentioned previously, in cycle 2 there are extensive supplemental sample buy-ins in Ontario, Manitoba and Alberta selected using Random Digit Dialing techniques. Therefore, in those three provinces, there exist samples from two frames, the regular (core) NPHS frame and the RDD frame. It is necessary to combine the weights coming from the core and RDD units so that the resulting estimates do not doubly estimate the population. A dual-frame technique, an adaptation of the Skinner and Rao method (1996)<sup>5</sup> is used for this purpose. The details of this adjustment are described in Section 11.4.

**Adjustment 6A: Post-stratification Adjustment**

This adjustment ensures that the final weights sum to the 1996-97 population totals at the province/age/sex level. The age-sex categories are people aged 0-11, 12-24, 25-44, 45-64, and 65 and older—males and females. This adjustment is given by

$$\frac{\text{population projection in a province / age / sex category}}{\text{sum of weights of respondent household members in a province / age / sex category}}$$

Again, it should be noted that households entirely composed of new immigrants since 1994 are not covered by the 1996-97 NPHS core sample, but are included in the population projections. Therefore these immigrants are implicitly treated as if they had the characteristics similar to the rest of the population.

Adjustment 7A: Noise Factor

For confidentiality reasons, a “noise” factor has been added to the weights of persons within the same household. This factor follows a uniform distribution and is added in such a way that the sum of the weights at the household level is respected.

The final household member weight on the general file is calculated as the “stripped” weight multiplied by all of the weight adjustments mentioned above.

**11.1.3 Weight Adjustments for Selected Members**

The final selected member weight on the health file is calculated as the “stripped” weight multiplied by the following weight adjustments. (Please note that these weight adjustments follow the same numbering scheme as in the previous section, for ease of comparison.)

Adjustment 1B: Household and Selected Member Non-response Adjustment

This adjustment compensates for complete household non-response, and for selected individuals within responding households who do not answer the selected members’ questionnaire. Similar to Adjustment 1A, weighting classes are constructed using information available on the selected member from 1994-95. Again, separate sets of weighting adjustment cells are created for each province but the input used for constructing the classes is slightly different. These classes are based upon data from all longitudinal respondents, not simply those who had responded to at least the general component in both years, as is used for the household member non-response adjustment. The adjustment is given by

$$\frac{\text{sum of weights for responding and non \& responding households in weighting class}}{\text{sum of weights for responding households in weighting class}}$$

Adjustment 2B: Interprovincial Migrations

This adjustment is made in the same manner as adjustment 2A. In cases where one extreme weight is generated within the provincial move pattern in question, the weights of all those selected members falling within the move pattern in question are reduced so that the sum of their weights equals the demographic projection of the number of movers within the move pattern in the past two years. If the adjustment would have resulted in an increase, the adjustment is not made, since that would augment the extreme weight even further.

Adjustment 5B: Integration of RDD Buy-ins

As in adjustment 5A, this adjustment is performed to combine the weights coming from the core and RDD units so that the resulting estimates do not doubly estimate the population. See Section 11.2 for more details.

Adjustment 6B: Post-stratification

This adjustment ensures that the final weights on the health file sum to the 1996-97 population totals at the province/age/sex level. The age-sex categories are people aged 2-11, 12-24, 25-44, 45-64, and 65 and older—males and females. This adjustment is given by

$$\frac{\text{population projection in a province / age / sex category}}{\text{sum of weights of respondent individuals in a province / age / sex category}}$$

Since no new longitudinal members were selected in 1996-97, the first age group only covers 2-11-year-olds. Again, it should be noted that households entirely composed of new immigrants since 1994-95 are not covered by the 1996-97 NPHS core sample, but are included in the population projections. Therefore these immigrants are implicitly treated as if they had the characteristics similar to the rest of the population.



## **11.2 Cross-sectional Weighting for the 1996-97 NPHS—Provinces with RDD Supplemental Samples**

### *Introduction*

To tap into the full cross-sectional potential of the data it is necessary to combine the core NPHS and supplementary RDD data into one large data set and have weights that reflect this combination of data sources. In some RDD provinces, extra questions are asked of the RDD respondents, or else some questions are asked of different age groups than for the core survey. The differences in covered populations for these questions mean that extra weights must be produced. In total there are three cross-sectional weights, one for analysis of the general file and two for analysis of the health file. The *Standard Household Member Final Weight* (WT56) is used when analyzing the data from the general form.

The *Standard Selected Adult Member Final Weight* and the *Standard Selected Child Member Final Weight* have been combined into WT66. This weight is used when analyzing most of the data for the health file for people aged 12 and older, and for the selected children aged 0-11 in Manitoba and Alberta. The second weight on the health file, WT66\_N, is applicable to two different populations, and therefore care must be taken when using this weight. One objective of this weight (the *HPS Selected Adult Member Final Weight*) is the analysis of health form data involving any of the questions that were part of the Health Promotion Survey. This applies for people aged 12 and older. The second purpose of this weight is the analysis of the children's health services questions (SVB6\_1 to SVB6\_5). These questions were only asked of the selected child in the Alberta and Manitoba RDD samples.

See Section 12.1 - Use of Weights for more details on when to use the various weights.

### *Overview*

The production of weights in RDD provinces requires four general steps. First the core weights up to but not including the post-stratification step are generated (see the previous section for details). Next the RDD weights are defined up to a similar stage. The RDD calculations are reviewed in this section. Next the two sets of weights are integrated to produce one set of weights using a method developed by Skinner and Rao (1996)<sup>5</sup>. Finally these weights are poststratified to population totals at the health area/age/sex level.

*The Skinner-Rao Integration Method*

This method is used to combine two independent and individually-weighted samples, such as the NPHS core and RDD samples, in order to come up with a single sample using a single set of composite weights. A factor  $a$  ( $0 < a < 1$ ) is determined so as to give the relative importance of each sample in the combined datafile. The core weight is multiplied by  $a$  while the RDD weight is multiplied by  $1-a$ . This approach applies in the case when the same population is being covered by both the core and RDD. However, there are cases in the NPHS where this is not true. One example is 0- and 1-year-old children who were eligible to receive a more detailed questionnaire in the RDD component but not in the core. In this case a more advanced use of the method is required.

The  $a$  value is not determined at the provincial level, but rather at what is called the collapsed health-area level. The provinces are divided into a number of large geographic regions, often separated between major city and outlying areas. In Alberta there are three collapsed health areas—Calgary, Edmonton and the rest of the province. In Manitoba there are two—Winnipeg and the rest of the province. In Ontario there are six areas—Toronto, Ottawa-Carleton and four other geographic regions. Collapsed areas are used so that their sample sizes become large enough to give stable input values for the calculation of  $a$  while still allowing for differences that may occur between distinct areas (large cities vs. the rest of the province). For more information on the method itself, see Skinner and Rao (1996)<sup>5</sup>.

*The Calculation of RDD Weights*

The weighting procedures used to calculate the RDD weights are similar to those used in the core. An initial RDD weight is generated based on the probability of selecting that telephone number. Several weight adjustments are then made to account for those selected numbers that were not sent to the field, or to account for other occurrences during data collection, such as households with multiple telephone lines. In many cases the adjustments are defined in terms of weighted ratios, with the weights created in the previous step being used to determine the magnitude of the adjustment. The resulting weight is then the previous weight multiplied by the weight adjustment.

**11.2.1 RDD Basic Weights**

Random Digit Dialing was implemented using the Elimination of Non-Working Banks Method, which gives rise to a stratified simple random sample (without replacement) of residential telephone lines. The entire province was divided into several RDD strata (see Chapter 5 - Sample Design for more details), which were not related to the strata used in the core design in the province. Each month a sample of telephone numbers was selected from the RDD frame. An initial

monthly weight is given by the inverse probability of selecting the particular telephone number from the list of working banks of numbers during that month. These monthly weights are converted to overall basic weights by multiplying at the stratum level the weights from month  $i$  by

$$\frac{\text{stratum sample size from month } i}{\text{total stratum sample size}}$$

If there was no change in the sampling rates or in the frame (the list of working banks) throughout the life of the survey then all of the weights within the stratum will be equal. If there has been a change, then the increase or decrease in the sampling rate or number of working banks will be reflected by higher or lower weights for the records from the month in question.

### **11.2.2 Further Weight Adjustments to the Basic Weights**

#### Adjustment 1: Manitoba Subsampling

In each stratum a theoretical “hit rate” is defined. This measures the expected number of residential telephone numbers in the stratum (as determined from a telephone billings file) compared with the total possible number of telephone numbers from the stratum after the Elimination of Non-Working Banks Method has been implemented. In several strata in Manitoba the hit rate was very low. This meant that a large sample would be drawn and most of the numbers would be nonworking, resulting in a very inefficient sample. In these regions, it was decided that approximately one-half of the sample drawn that did not appear as a residential number on the billings file would be discarded. (Eliminating all of them would exclude new telephone numbers and unlisted numbers from the sample, presenting a potential bias problem.) Thus, a multiplicative factor of

$$\frac{\text{total sample of non \& residential phone numbers in stratum}}{\text{sample of non \& residential phone numbers sent to the field in stratum}}$$

is applied to the non-residential numbers from the stratum that were sent to the field for interviewing. This factor is always close to 2.

Adjustment 2: Duplicate Telephone Numbers with the Core Sample

Since the core and RDD surveys covered the same population, but the samples were drawn independently, it was possible that a household could be selected in both components. To avoid this as much as possible, the selected RDD telephone numbers were compared with those on the core sample database to check if there were any matching numbers. If a match was found, the RDD number was not sent into the field. These households have to be accounted for in the weights of those numbers that were sent to the field. Those numbers have their weights adjusted by a factor of

$$\frac{\text{sum of weights of all households}}{\text{sum of weights of records sent into the field}}$$

This adjustment is done at the stratum level.

Adjustment 3: Household Non-response

To adjust for entire households that did not respond, the following adjustment is made:

$$\frac{\text{sum of weights for responding and non & responding households}}{\text{sum of weights for responding households}}$$

Again, this adjustment is performed at the stratum level. If the adjustment factor is high (greater than 2.5) then the offending stratum is collapsed with neighbouring strata until the adjustment falls below the 2.5 threshold.

Adjustment 4: Households Without Telephones

Obviously people without telephones in their dwelling have no chance of being selected and interviewed in the RDD survey. By simply ignoring this factor, a bias against people with characteristics similar to the non-telephone population will emerge. Ciok (1993)<sup>1</sup> has done a study on the non-telephone population in Canada. He divided the population into five household types and examined the non-telephone population in each category. The five groups are

- i) one person households—residents 65 years old or older
- ii) two or more person households—all residents aged 65 or older
- iii) one person households—residents less than 65 years of age
- iv) single parent households
- v) other households.

Using this information as well as the provincial non-telephone household percentages, adjustment factors were produced for each province within the five groups identified by Ciok. The weights are then multiplied by the appropriate non-telephone adjustment factor.

#### Adjustment 5: Multiple Telephones

Dwellings with more than one residential telephone line have a greater chance of being selected in the RDD sample. The total number of residential phone numbers in the household is collected during the interview. The weight resulting from the previous adjustments is then multiplied by the inverse of the number of residential lines. Thus, a household with two residential lines will have its weight divided by 2.

The weights resulting from this adjustment are called the “**RDD demographic weights**”.

### **11.2.3 Further Weight Adjustments for Household Members**

In calculating the RDD demographic weight, all of the adjustments have been made for household- or dwelling-level characteristics. Starting at this point, weight adjustments are made for individuals instead. The adjustments in this section relate to the household members’ weights, that is, the weights used when working with the data from the general component that was administered to all household members.

#### Adjustment 1H: Collective Dwellings

Collective dwellings, in which 10 or more unrelated persons live together, are in scope for the core survey but out of scope for the RDD component. Since the two pieces are combined at a later point, it is important that the weights for the RDD portion reflect the collective dwelling population as well. As with the non-telephone population, if a single adjustment factor is used over the entire population it will result in some bias.

In an effort to get an improved adjustment factor, information on collective dwelling residents was retrieved from the 1991 Census. The population was divided into 10 age-sex categories (males and females, aged 0-11, 12-24, 25-44, 45-64 and 65 and older) and three marital categories (married, single, divorced/separated/widowed). Within each of these categories, the percentage of people living in a collective dwelling was calculated. Note that some dwellings that the Census defines as a collective such as hospitals, jails and military bases

are not in scope for the NPHS. Only those categories common to the Census and NPHS were used in the calculation of the number of collective dwelling residents. Within each province, RDD demographic weights for each age/sex/marital status category were adjusted by a factor of

*( 1 % rate of people in age / sex / marital status category in collective dwellings ) .*

#### **Adjustment 2H: Household Member Non-response**

This adjustment compensates for individuals within responding households who do not respond to the household member questionnaire. The adjustment is equal to

$$\frac{\text{sum of weights of respondent and non \& respondent individuals in an age \& sex category}}{\text{sum of weights of respondent individuals in an age \& sex category}}$$

The age-sex categories are people aged 0-11, 12-24, 25-44, 45-64 and 65 and older—males and females. The adjustments are made within each health area. After this adjustment the resulting weight is called the “**RDD Household Member Weight**”.

#### **Integration of Core and RDD Household Members**

##### **The Standard Weight**

This weight is used when analyzing the questions on the general component. Since both the core and RDD cover all age ranges, a simplified version of the Skinner-Rao method is utilised to combine the Core Household Member Weight and the RDD Household Member Weight. The resulting weight is known as the “**Integrated Standard Household Member Weight**”. The final step in the creation of the Standard Weight is post-stratification.

#### **Adjustment 3HS: Post-stratification**

Monthly population projections are available for a number of age-sex categories at several different geographic levels. They are based upon the most recent Census as well as estimates of births, deaths, immigration and emigration. These data were used as an input to determine the totals at the health area/age/sex level within the three buy-in provinces.

The weights are then adjusted by a factor of

$$\frac{\text{population projection in a health area / age / sex category}}{\text{sum of weights of respondent household members in a health area / age / sex category}}$$

The age-sex categories are the same as those used in the member non-response adjustments, that is, people aged 0-11, 12-24, 25-44, 45-64 and 65 and older—males and females. The resulting weight is called the “**Standard Household Member Final Weight**” (WT56).

#### **11.2.4 Further Weight Adjustments for Selected Members**

In calculating the RDD demographic weight, all of the adjustments have been made for household- or dwelling-level characteristics. Starting at this point weight adjustments are made for individuals instead. The adjustments in this section relate to the selected members’ weights, that is, the weights used when working with the data from the health component that is administered only to the selected non-child, aged 12 and older, or to the selected child in the case of the Alberta and Manitoba RDD samples (recall that in these provinces both a non-child and a child are selected).

##### Adjustment 1S: Probability of Selection within the Household

Each member of the RDD responding household aged 12 and older had an equal chance of being chosen as the selected non-child. In addition, in Alberta and Manitoba, every child aged less than 12 within a household had an equal probability of selection (children could not be selected in Ontario). The weights of the selected individuals have to be increased to account for this extra stage of selection. This multiplicative adjustment is equivalent to the inverse probability of selection. For selected non-children this means the adjustment is equal to the number of non-children in the household, while for children it is equal to the number of children in the household.

##### Adjustment 2S: Collective Dwellings

This adjustment is done in the same way as adjustment 1H, except that the selected members weights are used this time.

Adjustment 3S: Selected Member Non-response

This adjustment compensates for selected individuals within responding households who do not answer the selected member questionnaire. The adjustment is equal to

$$\frac{\text{sum of weights of respondent and non-respondent individuals in an age \& sex category}}{\text{sum of weights of respondent individuals in an age-sex category}}$$

The age-sex categories are people aged 0-11, 12-24, 25-44, 45-64 and 65 and older—males and females. The adjustments are made within each health area. After this adjustment the resulting weight is called the “**RDD Selected Member Weight**”.

**Integration of Core and RDD Selected Members**

Two weights are required for selected members. For the non-children there is a standard weight used in most analyses and a Health Promotion Survey (HPS) weight used for certain questions that were not asked of the Alberta RDD sample. For the non buy-in provinces, the standard and HPS weights are identical. For children there is also a standard weight used in most cases as well as a health services weight used for certain questions on health services available to children that were only asked of the Alberta and Manitoba RDD samples.

As mentioned previously, the Standard Selected Non-child and Child Member Final Weights have been combined into one field on the microdata file, WT66. The HPS Selected Non-child Member Final Weight and the Health Services Selected Child Member Final Weight have been combined into WT66\_N on the microdata file. In all cases, a version of the Skinner-Rao integration method is used to combine the Core Selected Member Weight and the RDD Selected Member Weight. Finally a post-stratification step takes place to determine the final weight.

**The Non-child Standard Weight**

This weight is used when analyzing most of the questions on the health component. Since both the core and RDD sample cover the age ranges of 12 and older, the simplified version of the Skinner-Rao method is used to come up with a set of weights known as the “**Integrated Standard Selected Non-child Member Weight**”.



### Adjustment 3SSA: Post-stratification

Monthly population projections are available for a number of age-sex categories at several different geographic levels. They are based upon the most recent Census as well as estimates of births, deaths, immigration and emigration. This data was used as an input to come up with population totals at the health area/age/sex level within the three buy-in provinces. The weights are then adjusted by a factor of

$$\frac{\text{population projection in a health area / age / sex category}}{\text{sum of weights of respondent household members in a health area / age / sex category}}$$

The age-sex categories are the same as those used in the member non-response adjustments, that is, people aged 12-24, 25-44, 45-64 and 65 and older—males and females. The resulting weight is called the “**Standard Selected Non-child Member Final Weight**” (WT66).

### **The HPS Weight**

This weight is used when analyzing certain questions that are part of the Health Promotion Survey. The province of Alberta chose not to ask these questions of the RDD buy-in sample. For this reason there is no integration of core and RDD in Alberta. The Core Selected Member Weight is poststratified in the same manner as other non-RDD provinces (see Section 11.3.3). It is not poststratified at the health-area level because of the small sample size. In Manitoba and Ontario, the weight is the same as the Standard Household Member Final Weight. This weight is known as the “**HPS Selected Non-child Member Final Weight**” (WT66\_N).

### **The Child Standard Weight**

This weight is used when analyzing most of the extra questions administered to the selected child. Since there are no children aged 0 or 1 in the core, this is a special case of the Skinner-Rao integration method. In these two provinces, respondents from both surveys aged 2 to 11 are integrated together. For those aged 0 and 1, only the RDD records are available. The method also has a population count that ensures that the representativity of the two subgroups (ages 0, 1 and 2 to 11) is retained. The resulting weight is known as the “**Integrated Standard Selected Child Member Weight**”. It is not required in Ontario since children were not selected in the RDD sample. For this reason, the weights for Ontario can be calculated in the same manner as any other non-RDD

province, and are likewise poststratified at the provincial, not health-area, level.

Adjustment 3SSC: Post-stratification

Monthly population projections are available for a number of age-sex categories at several different geographic levels. They are based upon the most recent Census as well as estimates of births, deaths, immigration and emigration. This data was used as an input to come up with population totals at the health area/age/sex level within the three buy-in provinces. The weights are then adjusted by a factor of

$$\frac{\text{population projection in a health area / age / sex category}}{\text{sum of weights of respondent household members in a health area / age / sex category}}$$

The age-sex categories are the same as those used in the member non-response adjustments, that is, males and females aged 0-11. The resulting weight is called the “**Standard Selected Child Member Final Weight**” (WT66).

**The Child Health Services Weight**

This weight is used when analyzing questions on child health services. These questions were only asked of children selected in the Alberta and Manitoba RDD samples. There is no core data with which to integrate them. Outside of these two provinces, this weight is zero.

Adjustment 3SHC: Post-stratification

Monthly population projections are available for a number of age-sex categories at several different geographic levels. They are based upon the most recent Census as well as estimates of births, deaths, immigration and emigration. This data was used as an input to come up with population totals at the health area/age/sex level within the three buy-in provinces. The weights are then adjusted by a factor of

$$\frac{\text{population projection in a health area / age / sex category}}{\text{sum of weights of respondent household members in a health area / age / sex category}}$$

The age-sex categories are the same as those used in the member non-response adjustments, that is, males and females aged 0-11. The resulting weight is called the “**Health Services Selected Child Member Final Weight**” (WT66\_N).

### **11.3 1994-95-based Weighting Procedures for the Provinces Outside Québec**

To begin, the basic LFS weighting procedure is described below. A description of the other multiplicative weight adjustments necessary in the formation of the “stripped” weights follows.

#### **11.3.1 LFS Basic Weights**

The LFS uses a stratified multi-stage design (mainly two-stage, but in some cases, three-stage). In both cases, a sample of clusters is selected in each stratum using one of several probability proportional to size (PPS) sampling schemes. An LFS "cluster weight" is then calculated as the inverse probability of selecting a cluster, in accordance with the sample selection scheme. At the next (last) stage, dwellings are selected within sampled clusters using systematic sampling. A "dwelling weight" is calculated as the inverse probability of selecting a dwelling given that the cluster containing it is selected. An "LFS basic weight" is then given by the product of the cluster weight and the dwelling weight.

#### **11.3.2 Further Weight Adjustments to the Basic Weights**

##### Adjustment 1: Rotation Group Weight Adjustment

The full LFS sample is comprised of six "rotation groups". The NPHS requests sample from the LFS in terms of integral numbers of rotation groups (between 1 and 6), although a fractional number may actually be required to fulfil sample size needs (see Adjustment 3 below). Thus, the first multiplicative weight adjustment, which compensates for the integral number requested, is given by the following:

$$\frac{\text{number of rotations in an LFS stratum used by LFS (usually 6)}}{\text{integral number of rotations in an LFS stratum requested by NPHS}}$$

Adjustment 2: Cluster Growth Weight Adjustment

There may be clusters that experience growth between the time when a Census enumeration of the cluster takes place and the time when the cluster is listed for the LFS. The cluster selection probability is based on the Census enumeration figure, which may be out of date. This has the effect that the number of dwellings in the LFS sample increases very slightly with moderate growth in the housing stock. In clusters where substantial growth has taken place, subsampling is used to keep interviewer assignments manageable. The NPHS also institutes a similar subsampling of clusters that have experienced moderate growth. Thus, the second multiplicative weight adjustment is given by the inverse of this subsampling ratio in clusters where subsampling has occurred for either the LFS or the NPHS.

Adjustment 3: Stabilization Weight Adjustment

Stabilization is a means of capping the sample size within a stabilization area to prevent the associated costs from becoming too prohibitive. A "stabilization area" consists of clusters in the high-income and apartment frame, and of groups of strata in the regular frame. "Stabilization" addresses the problem of growth that occurs within a stabilization area. The growth is large enough to be a concern even after cluster growth adjustment, although no single cluster contributes substantially enough to the growth to be considered the root of the problem. This problem is remedied through subsampling within the stabilization area. In addition to regular stabilization, it is at this point that the fractional part of a rotation requested of the LFS, but not required by the NPHS, is also "stabilized out" through subsampling (see Rotation Group Weight Adjustment). Thus, the third multiplicative weight adjustment is given by the following:

$$\frac{\text{number of dwellings selected by the LFS within a cluster}}{\text{number of dwellings actually used by the NPHS within a cluster}}$$

Adjustment 4: Multiples Weight Adjustment

It sometimes happens that an interviewer discovers that a listing entry thought to constitute single private occupied dwelling in fact constitutes two or more private occupied dwellings. This may happen, for instance, when a basement apartment is attached to a dwelling but has its own separate entrance. In this case, since interviewing takes place in only one of the private occupied dwellings (selected at random), the weight associated with that dwelling is boosted. Thus, the fourth multiplicative weight adjustment is given by the number of private

occupied dwellings that the listing entry in question actually constitutes. For most listings, this adjustment factor will be one.

Adjustment 5: Household Non-response Weight Adjustment

Despite all the attempts made by the interviewers, some non-response at the household level is inevitable. Non-response encompasses any of the following situations: refusal, special circumstance, language barrier, no one at home, temporarily absent or computer problem. Non-response is compensated for by proportionally adjusting the weights of responding households.

It is also at this step that the removal of 1994-95 cross-sectional buy-in units, which were not to be surveyed in 1996-97, was accounted for. Households selected in 1994-95 as part of the buy-in are removed from both the numerator and denominator of the following calculations. This fifth adjustment is given by the following:

$$\frac{\text{sum of weights for sampled households in NPHS stratum / season combination}}{\text{sum of weights for respondent households in NPHS stratum / season combination}}$$

Note that this adjustment is made at the NPHS stratum level for each "season", i.e., summer and winter. Here, NPHS strata are groups of LFS strata. The adjustment was made at this level since it was the smallest geographic level that ensured stability. The adjustment was calculated separately for each season since the non-response rate was significantly different for each season. The "weights" referred to above are the LFS basic weight multiplied by all the adjustments to this point (i.e., weight adjustments 1 through 4). The adjustment is based on the assumption that the households that were actually interviewed represent the characteristics of those that should have been interviewed in each stratum/season combination.

Adjustment 6: Rejective Method Weight Adjustment

As discussed in Chapter 5 - Sample Design, in the last two quarters of data collection in 1994-95 a portion of the sampled households was screened out or rejected from the sample after determining that there were no youths or children residing within (i.e., no one less than the age of 25). These "rejected" households come from that portion of the "Children Sample" that are "screened" for household composition.

This methodology was implemented to compensate for an overrepresentation in the sample of members of small sized households and an underrepresentation of members of large sized households. The latter type of household often consists of parents and their children while the former type tends to consist of single people, older people or couples without children. Since some households containing no youths or children are screened out or "rejected", representation in the sample of households of this type comes solely from the "Adult Sample" and from the non-screened portion of the "Children Sample". Thus, to compensate for the "rejected" part of the sample, the weights for those households containing no youths or children from the "Adult Sample" and from the non-screened portion of the "Children Sample" are boosted by another multiplicative weight adjustment.

This sixth adjustment is given by the inverse of one minus the overall screening rate within a stratum. Note that in P.E.I., this adjustment was implemented a little differently since, among other reasons, the rejective method was applied in all four quarters of data collection rather than in the last two quarters only. Also note that this adjustment was not applied in apartment strata, high income strata and remote strata, since the rejective method was not implemented there.

### **11.3.3 Further Weight Adjustments for Selected Members**

Data from the selected members' questionnaire is obtained for only one member of a sampled household. The "stripped" weight for each individual is obtained as follows. First, the LFS basic weight is multiplied by weight adjustments 1 through 6, as well as weight adjustments 7B, 8B and 9B given below. If the selected person is a child less than 12 years of age who lives in a "Children" sample dwelling (see Chapter 5 - Sample Design for the definition of "Children" sample dwelling) then all of the children in the household to a maximum of four were administered the NLSCY questionnaire in 1994-95. Otherwise, the selected member was asked an additional set of NPHS questions. Several adjustments have to be made to account for this design and the non-response to this questionnaire.

#### Adjustment 7B: NLSCY Integration Weight Adjustment

In the last two quarters of data collection in 1994-95, the NPHS selected respondents for both the NPHS and NLSCY selected-member questionnaires. In sampled "Children" households that had children, all children up to a maximum of four (aged less than 12) were selected and administered the NLSCY questionnaire. One child was identified as the NPHS panel member for that household in future cycles. The children's data did not reside on the 1994-95

NPHS microdata file. For more details on integration with the NLSCY, see Chapter 5 - Sample Design. In 1996-97, all selected members including the panel children aged less than 12 in 1994-95 are surveyed by the NPHS, not the NLSCY.

In 1994-95, households containing children aged 12 or older were selected from the "Adult Sample" only. To compensate for the fact that households containing children coming from the "Children Sample" did not contribute to the estimates for selected individuals in 1994-95 but will in 1996-97, the weights for those households containing children sampled in the last two quarters that come from the "Adult Sample" had a special weight adjustment applied. This adjustment is given by the inverse of the proportion of the total sample assigned to the "Adult Sample". For those individuals aged more than 12, one adjustment is made at the cluster level. On the other hand, for those aged 12, a separate adjustment is made for groups of LFS strata (which usually correspond to NPHS strata), to be consistent with Adjustment 9B, which is also made at this level.

#### Adjustment 8B: Selected Member Inverse Selection Probability

As mentioned above, one member from each sampled household is chosen as the selected member. A weight adjustment must be made to reflect the selection and is given by the inverse selection probability. The original intention was that each member would be selected with equal probability given by the inverse of the number of members in the household. However, due to an error made in the CAI application, no 12-year-olds were selected in the first two quarters. To compensate, in the last two quarters, instead of each member of a household being selected with the same probability, 12-year-olds were given a larger probability of selection. In P.E.I, 12-year-olds were twice as likely to be selected as any other member aged 13 or more, and elsewhere in Canada, 1.75 times as likely to be selected as any other member aged 13 or more.

#### Adjustment 9B: 12-year-old Weight Adjustment

Due to the error mentioned above, 12-year-olds were selected only in the last two quarters of data collection. To obtain an accurate representation of 12-year-olds, their weights had to be adjusted to account for the first two quarters when they had no probability of being selected. This adjustment is made for groups of LFS strata, which usually correspond to NPHS strata, except for the cases of remote and high income strata. In households with children, 12-year-olds could be selected from the "Adult Sample" in all quarters, but were actually only selected from the "Adult Sample" in the last two quarters. Since, within most NPHS strata, 40% of the "Adult Sample" occurred in the last two quarters, the

weights of 12-year-olds selected in these two quarters were boosted by the inverse of this rate, or by 2.5.

On the other hand, in households with youths but no children, 12-year-olds could be selected from both the "Adult Sample" and the "Children Sample". However, in the first two quarters, they were not selected from the "Adult Sample" as they should have been due to the error mentioned above. Thus, in households with youths but no children, the weights of 12-year-olds were boosted by a multiplicative factor given by the ratio of the percentage of the total sample within an NPHS stratum where they should have been selected to the percentage of the total sample where they were actually selected, or by 1.6. Finally, in households with no youths or children, there were no 12-year-olds, so no adjustment was needed. Note that the rates differ somewhat in P.E.I., apartment strata, high income strata and remote strata.

The "weights" referred to above are the LFS basic weights multiplied by all the adjustments to this point (i.e., weight adjustments 1 through 6 as well as 7B, 8B and 9B). These are the 1996-97 "stripped" weights for the provinces outside Québec, and the starting point for the 1996-97 weighting.

#### **11.4 1994-95-based Weighting Procedures for Québec**

The National Population Health Survey used a subsample of the *Enquête sociale et de santé* (ESS) in its design (see Chapter 5 - Sample Design for more details). For this reason, the calculation of NPHS weights is tied to the weighting procedures used for the ESS. The following sections describe the ESS weighting procedures and the steps required to produce the 1996-97 "stripped" weights for NPHS members.

##### **11.4.1 ESS Weights**

The ESS contribution to the weights is calculated as follows:

###### ESS Cluster Weights

The ESS used a stratified multistage design. After several levels of stratification, clusters were selected from each stratum using probability proportional to size (PPS). The size measure used was the household count in the cluster based upon the 1986 Census. An "**ESS cluster weight**" can be calculated as the inverse probability of selecting a cluster.



### ESS Dwelling Weights

After selecting a cluster, a fixed number of dwellings were allocated to be selected from the cluster. Each dwelling in the cluster had an equal chance of being selected. The "**ESS dwelling weight**" is then the inverse of the probability of selecting the dwelling within the cluster multiplied by the ESS cluster weight.

#### **11.4.2 NPHS Basic Dwelling Weights**

There were two major steps to selecting the NPHS sample. First the subset of ESS clusters to be used in the NPHS had to be identified. Second the subset of ESS dwellings within each retained cluster had to be selected.

### Probability of Retaining an ESS Cluster for NPHS

As ESS strata were sometimes very small, NPHS strata were defined as comprising one or more ESS strata. A fixed number of clusters were allocated to be retained from each NPHS stratum. In cases where the NPHS stratum consisted of more than one ESS stratum, the allocation of clusters to ESS strata was proportional to the number of households in each ESS stratum in order to produce a PPS sample of clusters in each NPHS stratum. Fractional sample sizes were randomly rounded up or down to the next integer. Once the number of clusters to be retained from an ESS stratum had been determined, each cluster within the ESS stratum had the same probability of retention in most cases. The exceptions were clusters in which the number of dwellings grew by more than 150% between the 1986 Census and the 1992-93 ESS cluster listing. These clusters were given a higher probability of retention (either 100% or 40% greater probability of retention).

### Probability of Retaining an ESS Dwelling for NPHS

In clusters retained for the NPHS, only dwellings selected for the ESS were eligible to be selected for NPHS. Those dwellings that were out of scope for the ESS (businesses, collectives, demolished or abandoned) had a probability of one of being retained, in case they became in-scope for NPHS. From the ESS in scope dwellings, a fixed number of dwellings within each cluster were initially retained for the NPHS. A further subgroup of these selected dwellings was dropped because of their ESS household composition. The probabilities that a dwelling would be retained due to its household composition ranged from one-third for one-person households to 1 for households with children less than 12 years old.

The "**basic dwelling weight**" is the ESS dwelling weight times the inverse of the product of the ESS cluster retention probability and the ESS dwelling retention probability. The ESS dwelling retention probability includes both the probability of a dwelling being initially retained for NPHS and the probability of being retained due to its household composition.

### **11.4.3 Further Weight Adjustments to the Basic Weights**

#### Adjustment 1: Multiples Weight Adjustment

Sometimes when an interviewer visited a dwelling, he/she found an extra dwelling that was missed during cluster listing. An example of this might be a basement apartment. In this case each dwelling is known as a multiple. When this occurred, one dwelling was selected at random and interviewed. The weight of the selected dwelling is then adjusted by a multiplicative factor equal to the number of multiples.

#### Adjustment 2: Cluster Growth Weight Adjustment

In a few cases, clusters were relisted by NPHS. If there was a growth of 15-30% between ESS counts and NPHS counts, then a multiplicative weight adjustment of

$$\frac{NPHS \text{ count}}{ESS \text{ count}}$$

is made to each selected dwelling within the cluster. If the growth was less than 15% then the growth is assumed to be negligible and this adjustment is set to one. For all these dwellings, the multiples and cluster growth adjustments are multiplied by the basic dwelling weight to give a "**preliminary weight**".

If the growth was more than 30% then extra dwellings were selected for NPHS from the extra dwellings listed within the cluster. For these selected extra dwellings, the "**preliminary weight**" is the inverse of the product of the ESS cluster selection probability and NPHS cluster retention probability multiplied by

$$\frac{\text{number of extra dwellings listed}}{\text{number of extra dwellings selected}}$$

and the multiples adjustments. Since none of these dwellings was interviewed by the ESS, there is no way to categorize them into one of the ESS household composition categories.

### Adjustment 3: Household Non-response Weight Adjustment

To adjust for total nonresponding households, the following adjustment is made:

$$\frac{\text{sum of weights for respondent and non & respondent households}}{\text{sum of weights for respondent households}}$$

The weight here is the preliminary weight. A separate adjustment is done within a non-response weighting area. For the ESS in scope dwellings the non-response weighting areas are defined as an intersection of an NPHS stratum and ESS household type by quarter. For the dwellings that were added because the cluster had greater than 30% growth during NPHS relisting, the weighting area consists of the added dwellings within the cluster by quarter. The ESS out-of-scope dwellings are grouped into two non-response weighting areas by quarter for non-response adjustment purposes. The first group contains all dwellings with an ESS response code of 10 (demolished, vacant, abandoned). The second contains all dwellings with an ESS response code of 18 (collective or business). Multiplying the preliminary weight by the household non-response weight adjustment produces the "**demographic weight**".

#### **11.4.4 Further Weight Adjustments for Selected Members**

One member from each responding household is designated as the selected member. If this person is a child less than 12 years of age who lives in a "Children" sample dwelling (see Chapter 5 - Sample Design for the definition of "Children" sample dwelling) then all of the children in the household to a maximum of four were administered the NLSCY questionnaire in 1994-95. Otherwise, the selected member was asked an additional set of NPHS questions. Several adjustments have to be made to account for this design and the non-response to this questionnaire.

#### Adjustment 4: NLSCY Integration Weight Adjustment

In a "Children" sample household where a child is found, one child is chosen to be the selected member for the NPHS longitudinal panel. This children's data did not reside on the 1994-95 NPHS microdata file. An adjustment has to be made to account for the adults and youths in these dwellings who had no chance of being the selected member. This adjustment is only applied to adults and youths selected for the longitudinal panel in "Adult" dwellings where children were found by NPHS.

The adjustment is equal to the inverse subsampling rate for the "Adult" sample. The adjustment depends upon whether the ESS found children in the dwelling and upon the ESS urban density class to which the dwelling belongs. A separate adjustment is generated for dwellings where ESS found children and dwellings where ESS did not find children because the subsampling rate was different for these two categories. In the ESS Montréal and regional capitals classes, the adjustment is made at the cluster level, while in the ESS smaller urban agglomerations and rural sector classes, it is made at the NPHS stratum level. For an exception to this rule see "12-year-old Weight Adjustment" later in this section.

#### Adjustment 5: Selected Member Inverse Selection Probability

In a dwelling belonging to the "Children" sample in which there were no children less than the age of 12 or a dwelling belonging to the "Adult" sample, every member was originally intended to have an equal probability of being the selected member. However, due to a software error, 12-year-olds were not eligible to be selected in the first two quarters. To compensate for this they were given double the probability of being selected in quarters 3 and 4. A weight adjustment equal to the inverse probability of an individual within the household being the selected member is applied.

#### Adjustment 6: 12-year-olds Weight Adjustment

In order to get an accurate representation of 12-year-olds, their weight had to be increased to account for households where they were not eligible to be selected as a result of the software error. This adjustment is equal to the inverse probability that a 12-year-old was eligible to be selected from a dwelling where a person 12 or older was intended to be the selected respondent. Recall that in the Montréal and regional capitals classes, clusters are only covered in one quarter. In quarters 1 and 2, 12-year-olds were not eligible to be selected.

Therefore, in order for the weight adjustment to account for these ineligible 12-year-olds, it must be done at the NPHS stratum level rather than the cluster level. For consistency, both the integration and 12-year-old weight adjustment are calculated at the NPHS stratum level for 12-year-olds whatever the ESS class.

The 1996-97 "stripped" weight is calculated by multiplying the demographic weight by all of the adjustments made in this section.

**References:**

- <sup>1</sup> Ciok, R. (1993). The Non-telephone Population—A Preliminary Study. Internal Document. Statistics Canada.
- <sup>2</sup> Ernst, L. (1989). Weighting Issues for Household and Family Estimates. In *Panel Surveys*. (Eds. D. Kasprzyk, G. Duncan, G. Kalton and M.P. Singh.) New York: John Wiley and Sons, 135-159.
- <sup>3</sup> Kass, G.V. (1980). An Exploratory Technique for Investigating Large Quantities of Categorical Data. *Applied Statistics*, 29, 119-127.
- <sup>4</sup> Lavallée, P. (1995). Cross-sectional Weighting of Longitudinal Surveys of Individuals and Households Using the Weight Share Method. *Survey Methodology*, 21, 25-32.
- <sup>5</sup> Skinner, C.J. and Rao, J.N.K. (1996). Estimation in Dual Frame Surveys with Complex Designs. *Journal of the American Statistical Association*, 91, 433, 349-356.

## **12. File Usage**

This section starts with a discussion of the weight variables and explains how they should be used when doing tabulations on the public use microdata files. This is followed by an explanation of the variable naming convention that is employed for all cycles of the NPHS. The last part of the section discusses alternate approaches to data access available to analysts.

### **12.1 Use of Weights**

#### **12.1.1 Cross-sectional Weight - General File WT56**

Only one weight, WT56, appears on the general file. This weight is applicable to all age groups and provinces. This weight is based on the integrated core and RDD samples. ALL QUESTIONS ON THE GENERAL FILE SHOULD BE ANALYZED USING THIS WEIGHT.

(For a more detailed explanation on the creation of this weight, see sections 11.1 and 11.2 of the documentation on weighting.)

#### **12.1.2 Cross-sectional Weight - Health File WT66**

WT66 is applicable to all age groups and provinces, but its coverage of the youngest age group (0-11-year-olds) varies by province. In Manitoba and Alberta, RDD households that contained children had a child (0-11 years old) selected in addition to a non-child (12+ years old) to complete the health component. The coverage of WT66 applies to 0-11-year-old children. In the other provinces, only the children from the panel were followed in 1996-97, so the applicable age group is 2-11. This weight is based on the integrated core and RDD samples. WT66 IS THE PRINCIPAL WEIGHT VARIABLE ON THE HEALTH FILE AND SHOULD BE USED FOR ANALYSIS OF MOST VARIABLES. The exceptions follow in 12.1.3 where WT66\_N should be used.

(For a more detailed explanation on the creation of this weight, see sections 11.1 and 11.2 of the documentation on weighting.)

**12.1.3 Cross-sectional Weight - Health File WT66\_N**

Most of the questions on the health file can be analyzed using WT66. However, both the Health Promotion Survey (HPS) and Child Health Services questions, listed below, should be analyzed using WT66\_N instead.

Variables to analyze using WT66\_N:

<b>Health Promotion Survey Questions</b>
GHS6_11, GHS6_12, GHS6_13, GHS6_14, GHS6_15, GHS6_16A, GHS6_16B, GHS6_16C, GHS6_16D, GHS6_16E, GHS6_16F, GHS6_16G, GHS6_16H, GHS6_17, GHS6_18A, GHS6_18B, GHS6_18C, GHS6_18D, GHS6_18E, GHS6_18F, GHS6_18G, GHS6_18H, GHS6_18I
GHC6_21
HWS6_4, HWS6_5
HVS6_1, HVS6_2, HVS6_3, HVS6_4, HVS6_5, HVS6_6, HVS6_7, HVS6_8, HVS6_9
SMS6_8, SMS6_9, SMS6_13A, SMS6_13B, SMS6_13C, SMS6_13D, SMS6_13E, SMS6_13F, SMS6_13G, SMS6_13H, SMS6_14, SMS6_15, SMS6_16A, SMS6_16B, SMS6_16C, SMS6_16D, SMS6_17, SMS6_18A, SMS6_18B, SMS6_18C, SMS6_18D
ALS6_1, ALS6_2, ALS6_3, ALS6_4, ALS6_5, ALS6_6, ALS6_7
SHS6_4, SHS6_5, SHS6_6, SHS6_7, SHS6_7A
<b>Child Health Services Questions</b>
SVB6_1, SVB6_2, SVB6_3, SVB6_4, SVB6_5

The HPS questions were not asked of the RDD sampled units in Alberta and are applicable to non-children (12+) only. The Child Health Services questions were asked of the selected RDD children in Manitoba and Alberta only.

(For a more detailed explanation on the individual weight components used in the creation of this weight variable, see section 11.2.4 of the documentation on weighting, namely ‘The HPS Weight’, ‘The Child Standard Weight’, and ‘The Child Health Services Weight’.)

## **12.2 Variable Naming Convention**

In 1996-97, the NPHS adopted a variable naming convention which allows data users to easily use and refer to similar data from different collection periods and across survey components of the NPHS program. The following requirements were mandatory: restrict variable names to a maximum of 8 characters for ease of use by analytical software products; identify the survey occasion (1994-95, 1996-97, 1998-99...) in the name; and allow conceptually identical variables to be easily identifiable over survey occasions. For example, conceptually identical data on smoking were collected in 1994-95 and 1996-97. The variable names about smoking should only differ in the year position in the name that identifies the particular survey occasion in which they were collected. This convention will be followed throughout the longitudinal survey, and will be adopted by all NPHS surveys: the household survey, the institutional survey, the Northern survey, and supplements.

### **12.2.1 Variable Name Component Structure**

Each of the eight characters in a variable name contains information about the type of data contained in the variable.

Positions 1-2: Variable / Questionnaire section name  
Position 3: Survey type  
Position 4: Year/cycle variable appears  
Position 5: Variable type  
Positions 6-8: Variable number / name from questionnaire

For example:  
the variables DHC4\_AGE and DHC6\_AGE:

DH: in the Demographic and Household content section of the questionnaire;  
C: questions which are Core content on the household survey;  
4: appeared in 1994-95 cycle;  
6: appeared in 1996-97 cycle;  
\_: can be found on the questionnaire, and;  
AGE: the variable name.



**1996-97 NPHS PUBLIC USE MICRODATA DOCUMENTATION**

**12.2.2 Positions 1-2: Variable / Questionnaire Section Name**

The following values are used for the section name component of the survey:

AD	Alcohol dependence	IJ	Injuries
AL	Alcohol	IN	Income
AM	Administration (of the survey)	IS	Insurance
AP	Attitudes towards parents	LF	Labour force
BP	Blood pressure	MH	Mental health
CC	Chronic conditions	PA	Physical activities
CE	Contact exit	PC	Physical check-up
CI	Contact information (institutions, 1996-97)	PR	Province
CO	Coping (Alberta buy-in, 1994-95 and 1996-97)	PY	Psychological resources (self-esteem, mastery, sense of coherence)
DG	Drug use	RA	Restriction of activities
DH	Demographics and household	RP	Repetitive strain
DV	Dental visits	RS	Road safety
ED	Education	RT	Rationality (Manitoba buy-in, 1994-95)
ES	Emergency services	SD	Socio-demographics
EX	Eye examination	SH	Sexual health
FI	Balance and falling (institutions)	SM	Smoking
FS	Flu shots	SP	Sample identifiers (methodology)
GE	Geographic identifiers (methodology)	SS	Social support
GH	General health	ST	Stress
HC	Health care utilization	SV	Health services
HH	Household	TU	Tanning and UV exposure
HI	Health information	TW	Two-week disability
HS	Health status	VS	Violence / personal safety
HV	HIV	WH	Women's health: breast self-examination, breast examination, mammography and Pap smear
HW	Height and Weight	WT	Weights

**12.2.3 Position 3: Survey Type**

A	Asthma supplement
B	Province-specific buy-in content - children's questions
C	Core questions that will be repeated in each cycle
I	Institutions
K	Longitudinal children's questions
N	North (Yukon / NWT)
P	Province-specific buy-in content - adult questions
S	National supplement (Health Promotion Survey)
-	Cycle specific questions, not repeated in every cycle (stress in 1994-95, access to services in 1996-97)
3	Survey administration variables for household and demographic component (H03)
5	Survey administration variables for the General component (H05)
6	Survey administration variables for the Health component (H06)

**12.2.4 Position 4: Year / Cycle Variable**

4	1994-95
6	1996-97
8	1998-99
0	2000-01
2	2002-03
A	2004-05
B	2006-07
C	2008-09
D	2010-11
E	2012-13
F	2014-15

**12.2.5 Position 5: Variable Type**

_	Collected variable	A variable that appeared directly on the questionnaire
C	Coded variable	A variable coded from one or more collected variables (e.g., SIC, Standard Industrial Classification code)
D	Cross-sectional derived variable	A variable calculated from one or more collected or coded variables, usually calculated during head office processing (e.g., health status index)
F	Flag variable	A variable calculated from one or more collected variables (like a derived variable), but usually calculated by the computer application for later use during the interview (e.g., work flag)
G	Grouped variable	Collected, coded, suppressed or derived variables collapsed into groups (e.g., age groups)
L	Longitudinal derived variable	A variable calculated using variables from two or more survey cycles

**12.2.6 Positions 6-8: Variable Name**

In general, the last three positions follow the naming on the questionnaire. Numbers are used where possible: Q1 becomes 1. “Mark-all” questions use letters for each possible answer category: Q1 (mark all that apply) becomes 1A, 1B, 1C, etc. Demographic variables which are used frequently by analysts are identified by a three letter identifier, rather than by a question number; for example “age” is DHC6\_AGE in 1996-97. Where groups of questions with the same topic were collected in sections that had different section names on the questionnaire, position 6 is used to identify the subsection. For example, the first question on chronic stress was named ST\_4\_C1, the first question on childhood and adult stressors (traumas) was named ST\_4\_T1. Another example of this occurs in the general health questions for the Health Promotion Survey. These questions were separated into three sections for inclusion in the questionnaire

and the corresponding variable names reflect this, with position 6 indicating the section in which it appears.

### **12.3 Remote Access of Master Files**

Microdata files must meet stringent security and confidentiality standards required by the Statistics Act before they are released for public access in order to protect the confidentiality of the respondents participating in the survey. To ensure that these standards have been achieved, each microdata file goes through a formal review process to ensure that an individual cannot be identified. Rare values in variables that may lead to identification of an individual are suppressed on the file or are collapsed to broader categories so that individual disclosure is minimized. Frequently, these are the variables that are the most critical for doing a complete and comprehensive analysis of the survey data. Since a significant amount of resources is spent on collecting these data, ensuring that the microdata files reach their full analytical potential is important for a complete return on the statistical investment.

Remote access to the survey master file is one way to reap these benefits. Each purchaser of the microdata product can be supplied with a 'dummy' test master file and a corresponding record layout. With this, the user can spend time developing his or her own set of analytical computer programs using the test file to confirm that the routines are functioning correctly. At that point, the code for the custom tabulations should be sent via the Internet to [nphs@statcan.ca](mailto:nphs@statcan.ca). The code will be moved into Statistics Canada's internal secured network and processed against the appropriate master file of NPHS data. The results are screened for confidentiality and reliability concerns and, once these have been addressed, the output is returned to the client. There is no charge for this service.

A second approach for any client is the production of custom tabulations done by the Client Custom Services staff in Health Statistics Division. This service allows users who do not possess knowledge of tabulation software products to have access to the master file for the preparation of their own custom calculations. As with remote access, the results are screened for confidentiality and reliability concerns before release. Unlike remote access, there is a charge for this service.