

STATISTICS CANADA
NATIONAL POPULATION HEALTH SURVEY
HOUSEHOLD COMPONENT
CYCLE 1 to 7 (1994/1995 to 2006/2007)
LONGITUDINAL DOCUMENTATION

July 2008



Statistics
Canada

Statistique
Canada

Canada

NOTE TO USERS

This documentation accompanies the release of 7 cycles of longitudinal data for the Household Component of the National Population Health Survey (NPHS).

This document provides a wide range of information on the NPHS: objectives, content, sample design, collection, processing, weighting procedures, data quality, tabulation's guidelines and data access. Chapters 7, 8 and 11 give more details on the various subsets of respondents and their associated sampling and bootstrap weights.

This document sometimes refers to a specific cycle of NPHS by using the years in which it occurred. For information, here is the list of NPHS cycles with their corresponding years:

Cycle 1 = 1994/1995
Cycle 2 = 1996/1997
Cycle 3 = 1998/1999
Cycle 4 = 2000/2001
Cycle 5 = 2002/2003
Cycle 6 = 2004/2005
Cycle 7 = 2006/2007

This document is also intended for the share file users i.e. the Provincial Ministries of Health, Health Canada and the Public Health Agency of Canada. The share file now includes two subsets (square and full – see section 7) which include respondents who agreed to share the information collected as part of the NPHS (all cycles) and the corresponding sampling and bootstrap weights. These subsets of respondents are also part of the master file. However, they are to be used by the share partners. The share file users should disregard references specific to other subsets of respondents.

**Quick Links Towards Statistics Canada WEB Pages
NPHS, Household Component – Longitudinal**

Documentation, NPHS Household Component – Longitudinal

- [Cycle 7 – Definitions and methods](#)
- [Cycle 7, Reference Documents](#) - On this page, you will find links towards:
 - Topical index
 - Longitudinal Documentation (PDF version of this document)
 - Derived Variables Documentation
 - Derived Variables List
 - Content Summary Cycle 1 to 8
 - Data Dictionary (rounded frequencies)
- [Definitions and Methods \(by cycle\)](#)

Questionnaires

- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 1 \(1994/1995\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 2 \(1996/1997\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 3 \(1998/1999\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 4 \(2000/2001\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 5 \(2002/2003\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 6 \(2004/2005\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 7 \(2006/2007\)](#)
- [NPHS Questionnaire, Household Component, Longitudinal, Cycle 8 \(2008/2009\)](#)
- [NPHS content summary, Household Component, Longitudinal, Cycles 1 to 7](#)

NPHS Internet Publication

- [Healthy Today, Healthy Tomorrow? Findings from the National Population Health Survey](#)

Table of Contents

	Page
1. Introduction.....	1
2. Background	2
3. Objectives.....	3
4. Survey Content.....	4
4.1 Content Selection Criteria	4
4.2 Cycle 7 (2006/2007) Changes to Existing Content	5
4.3 Previous Interviews Variables Used as Additional Information in Cycle 7	5
5. Sample Design	6
5.1 Cycle 1 (1994/1995) Sample Design	6
5.1.1 Sample Allocation.....	6
5.1.2 The Rejective Method.....	7
5.1.3 Sample Selection.....	7
5.1.4 Sample Design in Quebec.....	8
5.2 Longitudinal Sample.....	9
6. Data Collection	11
6.1 Questionnaire Design and Data Collection Method	11
6.2 Tests	11
6.3 Interviewing.....	11
6.4 Non-Response and Tracing	12
7. Data Processing	14
7.1 Editing	14
7.2 Coding.....	14
7.3 Derived and Grouped Variables.....	15
7.4 Estimation and Weighting	16
7.5 Definition of the Longitudinal Response Pattern (LONGPAT)	16
7.6 Definition of Full/Complete Response and Non-Response.....	17
7.7 Subsets of Respondents	17
8. Weighting	18
8.1 Longitudinal Weighting.....	18
8.1.1 Starting Point: Cycle 1 (1994/1995) Stripped Weights	19
8.1.2 Adjustments to Create the Different Longitudinal Weights	19
8.1.2.1 Longitudinal Square Weight (WT64LS)	19
8.1.2.2 Longitudinal Full Weight (WT6BLF).....	19
8.1.2.3 Longitudinal Square Share Weight (WT6BSLS).....	27
8.1.2.4 Longitudinal Full Share Weight (WT6BSLF).....	28
9. Data Quality.....	30
9.1 Sampling Errors	30
9.1.1 Bootstrap Method for Variance Estimation	30
9.2 Non-Sampling Errors.....	31
9.2.1 Response Rates	32
9.2.1.1 Cycle 1 (1994/1995) Response Rates.....	32
9.2.1.2 Cycle 2 (1996/1997) and Cycle 3 (1998/1999) Response Rates	33
9.2.1.3 Cycle 4 (2000/2001), Cycle 5 (2002/2003), Cycle 6 (2004/2005) and Cycle 7 (2006/2007) Response Rates.....	35

9.2.2	Refusal Rates	37
9.2.3	Unable to Trace Rates.....	38
9.2.4	Attrition Rates	38
9.2.4.1	Attrition Rates Based on the Longitudinal Full Subset.....	39
9.2.5	Item Refusal and “Don’t Know” Rates	41
9.2.5.1	Refusal and “Don’t Know” Rates by Item.....	41
9.2.5.2	Refusal and “Don’t Know” Rates by Respondent	42
10.	Guidelines for Tabulation, Analysis and Release	44
10.1	Rounding Guidelines	44
10.2	Sample Weighting Guidelines for Tabulation	44
10.2.1	Definitions of Types of Estimates: Categorical vs. Quantitative	45
10.2.2	Tabulation of Categorical Estimates.....	46
10.2.3	Tabulation of Quantitative Estimates.....	46
10.3	Guidelines for Statistical Analysis	46
10.4	Release Guidelines	47
11.	Using the Longitudinal Master File	49
11.1	Use of Longitudinal Weights	49
11.2	Ensuring the Reliability of Estimates with the Use of Bootstrap Weights.....	49
11.3	Variable Naming Convention	50
11.3.1	Variable Name Component Structure.....	50
11.3.2	Positions 1-2: Variable Name / Questionnaire Section Name	51
11.3.3	Position 3: Survey Type.....	53
11.3.4	Position 4: Cycle (years).....	53
11.3.5	Position 5: Variable Type.....	54
11.3.6	Positions 6-8: Variable Name	54
12.	Access to NPHS Data	55
12.1	Research Data Centres.....	55
12.2	Remote Access	55
12.3	Data Liberation Initiative.....	55
12.4	Analytical Reports and Tabulations.....	55
13.	An Analytical Technique for Longitudinal Survey Data.....	57
13.1	Cycle Twinning Approach	57
13.2	Creation of the Modified Subset.....	59
13.3	Methodological Aspects of the Cycle Twinning Approach	60
13.4	An Example of How to Use the Cycle Twinning Approach: Quitting Smoking	61

Tables List

Table 5.A	Longitudinal Sample Size by Province
Table 7.A	Distribution of deaths by year of death
Table 7.B	Subsets of Respondents
Table 8.A	Subsets of Respondents and Corresponding Sampling Weights and Flags
Table 8.B	Variables for Cycle 2 Non-response Adjustment
Table 8.C	Variables for Cycle 3 Non-response Adjustment
Table 8.D	Variables for Cycle 4 Non-response Adjustment
Table 8.E	Variables for Cycle 5 Non-response Adjustment
Table 8.F	Variables for Cycle 6 Non-response Adjustment
Table 8.G	Variables for Cycle 7 Non-response Adjustment
Table 8.H	Variables for the non-response adjustment
Table 9.A	Relevant information for calculation of response rates for Cycle 1
Table 9.B	Cycle 1 Response Rates
Table 9.C	Relevant information for calculation of response rates for Cycles 2 and 3
Table 9.D	Panel Response Rate for Cycles 2 and 3
Table 9.E	Relevant information for calculation of response rates for Cycles 4 to 7
Table 9.F	Panel Response Rate for Cycles 4, 5, 6 and 7
Table 9.G	Refusal Rates by Cycle
Table 9.H	Unable to Trace Rates by Cycle
Table 9.I	Attrition Type by cycle – Longitudinal Full subset of respondents
Table 9.J	Refusal and “Don’t Know” Rates by module
Table 9.K	Refusal and “Don’t Know” Rates by Respondent
Table 10.A	Sampling Variability Guidelines
Table 11.A	Subsets of Respondents and Corresponding Bootstrap Weights Files
Table 11.B	“Constant” Longitudinal Variables
Table 13.A	Example of profiles of response from fictitious NPHS respondents
Table 13.B	Attrition rates
Table 13.C	Example of the modified subset structure for twinning

Appendix List

Appendix A:	NPHS Household Component, Changes to the Questionnaire for Cycle 7 (2006/2007)
Appendix B:	NPHS Household Component, Examples of Cycle 7 (2006/2007) Data Feedback and Follow-up Questions

Other Reference Documents

Questionnaire
Record Layout
Alphabetic Index, Topical Index
Data Dictionary
NPHS Derived Variables Documentation, Cycles 1 to 7
NPHS, Derived Variables List, Cycles 1 to 7
NPHS, Content Summary, Cycles 1 to 8

1. Introduction

The National Population Health Survey (NPHS) is designed to collect "longitudinal" information on the health of the Canadian population and related socio-demographic information. The first cycle of data collection took place in 1994/1995. The survey will continue every second year thereafter for 10 cycles. The NPHS fulfilled both cross-sectional and longitudinal needs during its first three cycles, and then with Cycle 4 (2000/2001) the NPHS Household component became strictly a longitudinal survey. The cross-sectional component of the Population Health Surveys Program has been taken over by the Canadian Community Health Survey (CCHS).

The NPHS is now composed of only one component: the household component. Since 2000/2001, the North Component is conducted by CCHS rather than NPHS. After 5 cycles of collection (1994/1995 to 2002/2003), the Health institutions component was terminated due to the large number of deaths in the sample.

The target population of the NPHS Household component includes household residents in the ten Canadian provinces in 1994/1995 excluding persons living on Indian Reserves and Crown Lands, residents of health institutions, full-time members of the Canadian Forces Bases and some remote areas in Ontario and Quebec.

The Household component of NPHS has completed seven cycles: Cycle 1 (1994/1995), Cycle 2 (1996/1997), Cycle 3 (1998/1999), Cycle 4 (2000/2001), Cycle 5 (2002/2003), Cycle 6 (2004/2005) and Cycle 7 (2006/2007).

The Cycle 7 NPHS Household component collected in-depth information on the health of the longitudinal respondents who were randomly selected in Cycle 1 and demographic information about all members of the longitudinal respondents' household. The questionnaire includes questions related to health status, use of health services, determinants of health, chronic conditions and activity restrictions. Socio-demographic information is also collected; it includes age, sex, education, household income and labour force status.

This document has been produced to facilitate the use of Cycles 7 Longitudinal Master and Share Files from the Household component. **These files include data collected from respondents from Cycles 1 to 7.** The files are described in more detail in the following chapters.

Information requests, questions about data access (Research Data Centres, Remote Access), custom tabulations, general data support or about the data sets or their use should be directed to the:

Data Access and Information Services, Health Statistics Division:

Tel: 1-613-951-1746, E-mail: nphs-ensp@statcan.ca, Fax: 1-613-951-0792

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 systems and the commensurate requirement for information 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.

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. In addition, it was to be responsive to changing requirements, interests, and policies.

Since its beginning in 1994, the National Population Health Survey (NPHS) has been providing unique information on the health of Canadians by responding to the need for information on health dynamics. The NPHS is a longitudinal survey with a sample of 17,276 individuals spread out in the ten provinces across Canada. Every two years, these same individuals provide current and in-depth information on their physical and mental health status, use of health care services, physical activities, life at workplace and social environment. Over the years of follow-up, the data have shown how a wide range of factors can contribute to the improvement or deterioration of health.

Whereas data collected from people at a single point in time provides a snapshot, NPHS longitudinal data reveals the transitions towards good or poor health. The richness of NPHS' data is that it also allows evaluation of the relationships between socio-economic and demographic characteristics of individuals with their health status and its evolution over time.

3. Objectives

The objectives of the NPHS are to:

- aid in the development of public policy by providing measures of the level, trend and distribution of the population's health status;
- provide data for analytic studies that will assist in understanding the determinants of health;
- collect data on the economic, social, demographic, occupational and environmental correlates of health;
- increase the understanding of the relationship between health status and health care utilization, including alternative as well as traditional services;
- provide information on a panel of people who will be followed over time to reflect the dynamic process of health and illness;
- provide the provinces and territories and other clients with a health survey capacity that will permit supplementation of content or sample¹;
- allow the possibility of linking survey data to administrative data that are routinely collected, such as vital statistics, environmental measures, community variables, and health services utilization.

¹ Due to the longitudinal nature of NPHS the sample option is no longer available. CCHS is now providing this possibility.

4. Survey Content

The above noted objectives provided a broad direction for NPHS, particularly concerning the type of information to be collected. The first section of this chapter discusses the general criteria used for the selection of survey content and gives a broad summary of the questionnaire sections. The next section describes briefly the changes made to the NPHS content for Cycle 7. The last section provides information about variables from previous interviews that are used as additional information in Cycle 7.

Note that the NPHS content is quite stable since Cycle 6. Considering the longitudinal nature of the survey no new content was added since then. However, some previous modules or questions are repeated in the questionnaire from cycles to cycles in order for users to be able to do comparison over time.

NOTE: On page *iii*, you will find links to the NPHS questionnaires and to the NPHS content summary for seven cycles available on Statistics Canada Website..

4.1 Content Selection Criteria

The NPHS content was selected according to the following criteria:

- 1) Information should relate to and help monitor the health goals and objectives of the provinces. 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.²
- 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 Cycle 1, one person in each household was randomly selected as the longitudinal respondent to answer an in-depth questionnaire on health (Health component, H06). During the first three cycles, some information (demographic and health) was also collected about all member of the longitudinal respondent's household (General component - H05) and about the dwelling.

Starting in Cycle 4, the General and Health component questionnaires were combined into a unique questionnaire completed by the longitudinal respondents. Only the following demographic information: age, sex, marital status, relationships; is collected about the

² No new content was added to NPHS questionnaire since Cycle 6. This criteria is not valid anymore.

other members of the household. Information about the dwelling is also collected.

Reflecting the above criteria, the NPHS questionnaire includes questions related to health status, use of health services, determinants of health, chronic conditions and activity restrictions, and demographic and socio-economic status. For example, health status is measured through questions on self-perception of health, functional ability, chronic conditions, and activity restriction. The use of health services is measured through questions on visits to health care providers (traditional and non-traditional), hospital care and on use of drugs and other medications. Health determinants that are explored include smoking, alcohol use and physical activity. Questions are asked on preventive tests and examinations, which probed for frequency and reasons for use. Demographic and socio-economic information include age, sex, education, ethnicity, race, household income and labour force status.

4.2 Cycle 7 (2006/2007) Changes to Existing Content

As in previous cycle, questions on health are asked first and they are followed by the socio-economic questions (language, education, labour force status, and income).

The Cycle 7 focus content was incorporated in the most suitable place in the questionnaire. The Cycle 7 focus content includes Food Choices (was part of cycle 5), Childhood and Adult Stressors “trauma” (was part of cycle 4), and Food Insecurity (was part of cycle 3).

Appendix A details changes made to the Cycle 7 questionnaire.

4.3 Previous Interviews Variables Used as Additional Information in Cycle 7

In order to reduce respondent burden, questions to which the answer was already known and that would not change over time (e.g., country of birth) are not repeated.

Variables that could change over time if certain actions had occurred (e.g., level of education), were updated only if appropriate.

Some answers from earlier cycles were brought forward into the Cycle 7 interview. This proved to be a valuable tool resulting in better quality collected data. For instance, previous information on selected chronic conditions was recalled for the respondent in order to explain any change. For more information, please see Appendix B.

5. Sample Design

The target population of the NPHS Household component includes household residents in the ten Canadian provinces in 1994/1995 excluding persons living on Indian Reserves and Crown Lands, residents of health institutions, full-time members of the Canadian Forces Bases and some remote areas in Ontario and Quebec. This chapter describes the Cycle 1 sample design and explains how the sample of 17,276 persons was selected.

5.1 Cycle 1 (1994/1995) Sample Design

The Labour Force Survey (LFS) sample design, redesigned in 1991, was used as the basis for the sample design in all provinces except Quebec where the NPHS sample was selected from households already being interviewed by Santé Québec for the 1992-1993 *Enquête sociale et de santé* (ESS).

Three factors shaped the sample design of the household component sample:

- the targeted national and provincial sample sizes;
- the decision to select one member per household to make up the longitudinal panel;
- the choice of the LFS sample design 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 Quebec.

5.1.1 Sample Allocation

The NPHS initially had a target sample size of 19,600 respondent households. It was further agreed by 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)³. According to this scheme the sample was allocated proportionally to $\sqrt{(0.804W_h^2 + 1/12^2)}$, where W_h is the 1991 Census proportion of households in province h , $h=1, \dots, 10$. This allocation determined the base sample size for each province. Four provinces chose to increase their allotted sample size for the first cycle through the option of buy-in of additional units with increased funding, for cross-sectional purposes. These additional units were not retained for the longitudinal sample.

³ Kish, L. (1988). Multipurpose Sample Design, *Survey Methodology*, 14, 19-32.

5.1.2 The Rejective Method

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 under-represent persons coming from large households, typically parents and children, since they had less chance of being chosen and over-represent persons coming from small households, often single people or the elderly.

Thus, a rejective method was adopted to increase the representation of parents and youths in the panel. To do so, a portion of the sample was pre-identified for screening. After their member roster was completed, screened households that had no member less than 25 years of age were eligible for rejection 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 preliminary 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 method was also not applied in remote regions because of the high contact costs there.

5.1.3 Sample Selection

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

In all provinces except Quebec, 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 sample 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 sampled 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 a certain amount of non-response was expected, and a portion 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 a desired level, 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 this was in order to meet the specific sub-provincial sample size requirements for cross-sectional purposes in the first cycle. It was also required that the number of clusters selected per stratum be a multiple of four for variance estimation and seasonal representativeness (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. Once strata were grouped, their sample clusters were also grouped to form replicates.

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. Two design variables named "Stratum" and "Replicat" can be found on the Master file, where Stratum represents the LFS stratum, and Replicat represents the NPHS replicates.

5.1.4 Sample Design in Quebec

In Quebec, the NPHS sample was selected from dwellings participating in a Santé Québec health survey: the 1992/1993 *Enquête sociale et de santé* (ESS). The survey sampled 16,010 dwellings using a two-stage sample design similar to that of the LFS. The province was divided geographically by crossing 15 health areas with four urban density classes (Montreal Census Metropolitan Area, regional capitals, small urban agglomerations and the rural sector). In each area, clusters were stratified by socio-economic characteristics and were 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 Montreal the allocation was proportional to $\sqrt{(2W_h^2 + 1/4^2)}$, 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 National Longitudinal Survey of Children and Youth (NLSCY), affected the sub-sampling in Quebec 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.

5.2 Longitudinal Sample

The longitudinal sample, also called the longitudinal panel or simply the panel, is composed of the 17,276 persons that were selected in Cycle 1 and had completed at least the General component of the questionnaire in Cycle 1. It also includes 2,022 children from the first cycle (1994/1995) of the National Longitudinal Survey of Children and Youth (NLSCY). These children were interviewed by the NLSCY for the NPHS in Cycle 1 and are interviewed by the NPHS since the second cycle. This longitudinal sample (17,276), has been and will be surveyed in all NPHS cycles. Additional samples added to Cycles 1, 2 and 3 for cross-sectional purposes are not part of the longitudinal sample.

The longitudinal sample is not renewed over time. No panel members were or are to be classified out-of-scope. The longitudinal sample size remains the same (17,276) for all cycles. Consequently, for Cycle 7, all longitudinal panel members were 11 years old and over and the longitudinal sample did not contain anyone who has immigrated to Canada after 1994/1995.

The number of people answering the survey slightly decreases from one cycle to the next due to attrition caused by non-response (for example, refusals and individuals that were untraceable). Despite the attrition, the longitudinal sample is still representative of the

1994/1995 population. The attrition, being relatively small (see Section 9.2.4), should not lead to large increases in the variance of estimates. Note that panel members who died and panel members who moved to a health institution are still part of the longitudinal sample and are considered as respondents (see section 7.6). Therefore, these persons do not contribute to the attrition of the NPHS longitudinal panel.

Table 5.A presents the sample size of the longitudinal sample by province in 1994/1995. It also shows the number of people that provided a full response to all seven cycles of NPHS.

Table 5.A: Longitudinal Sample Size by Province

Province	Longitudinal Sample Cycle 1 (1994/1995)	Number of Respondents Full Response in Cycles 1 <u>to</u> 7
Newfoundland	1,082	746
Prince Edward Island	1,037	719
Nova Scotia	1,085	704
New Brunswick	1,125	728
Quebec	3,000	1,890
Ontario	4,307	2,546
Manitoba	1,205	805
Saskatchewan	1,168	824
Alberta	1,544	979
British Columbia	1,723	1,051
Total	17,276	10,992

6. Data Collection

6.1 Questionnaire Design and Data Collection Method

The survey questions were designed for computer-assisted interviewing (CAI), which means that, as the questions were developed, the associated logical flow into and out of the questions was specified, along with 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 is 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 advantage is automatic insertion of reference periods based on current dates. Pre-filling of text or data based on information gathered during the current interview or previous cycles' interviews allows the interviewer to proceed without having to search back for previous answers. This type of pre-fill 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 is customised to the respondent based on the data collected.

6.2 Tests

The CAI application was extensively tested in-house in order to identify any errors in the program flow and text. From Cycles 1 to 6 two field tests were conducted. The tests involved four of Statistics Canada's Regional Offices. The main objectives of the two tests were to observe respondent reaction to the survey, to obtain estimates of time for the various sections, to study response rates and to test feedback questions. Field operations and procedures, interviewer training, and the CAI application (i.e., the questionnaire on computer) were also tested. Beginning in Cycle 7 only one field test is conducted. The samples from the two field tests were combined. The majority of long-term non-respondents were removed from the sample. A small number was kept in order to be able to test tracing procedures among other things. The objectives of the test remain the same as before.

6.3 Interviewing

In Cycle 7, data collection for the household component was divided into four quarters (starting in June, August and October 2006, and January 2007). An additional collection period was held in June 2007 with further follow-up of non-respondents from previous quarters.

In Cycle 7, interviewers working in Statistics Canada Calling Centres located in Edmonton, Sturgeon Falls, Sherbrooke and Halifax performed data collection.

A special collection is conducted for the panel members residing in health care institutions. The interview for these respondents was conducted in person using a paper questionnaire. The NPHS health care component questionnaire was used and 272 respondents have been interviewed that way in Cycle 7. They are identified as such on the master and share

files (see section 7.5). However, the health institutions component questionnaire collects less information than the household component questionnaire. Missing variables for these respondents are coded to “6”, i.e. not applicable on the household component files.

All interviewers were employees hired and trained by Statistics Canada specifically to carry out surveys. NPHS data collection was performed under the supervisory and control structure put in place by Statistics Canada. All interviewers attended a training session that focused on NPHS content and they received an Interviewer’s manual for use as a reference tool.

Each living longitudinal panel member received by mail a letter announcing the start of NPHS Cycle 7 data collection. At the same time, the household component’s respondents received a brochure that presents general information about the survey as well as some results from the NPHS, some newspaper extracts (Breaking News) citing NPHS results and to thank them for their participation, a tool that allow to easily compute the body mass index (BMI). Furthermore, NPHS information on the Statistics Canada’s Website for the survey participants was also available.

In Cycles 6 and 7, a follow-up survey was conducted with a sub-sample of the November field test respondents. The objective of this follow-up was to verify the following elements with the respondents:

- Clearness and usefulness of the Brochure, Breaking News, BMI tool, Statistics Canada’s Website for survey participants)
- Impact of these documents on their motivation to participate
- Interest to receive additional NPHS information (Cycle 6)

These follow-up surveys showed that the majority of the respondents who received the documents were satisfied and read them. However, the impact on their participation is minimal. The BMI tool was modified based on received suggestions. Finally, the majority of consulted respondents wanted to receive additional survey information. Starting in Cycle 6, a question was added at the end of the questionnaire to identify these respondents. When the thank you material is sent, additional information is included for interested respondents.

In general, respondents from the household component are contacted by telephone. In fact, 99% of the interviews in Cycle 7 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, upon request by the respondent or if the respondent resided in a health care institution. The total interview time averaged just under an hour.

Information about all household members (age, sex, and relationships between members) was obtained from the longitudinal respondent. Proxy reporting for the longitudinal respondent aged 12 and over was allowed only for reasons of illness or incapacity. Such proxy reporting accounted for 5.7% of the information collected for respondents aged 12 years and older. On the other hand, all interviews for respondents under 12 years old were done by proxy.

6.4 Non-Response and Tracing

Many strategies were put in place to reduce the number of non-response cases. For example, the maximum assignment size for an interviewer was set to avoid overburdening

interviewers and was based on the experience from previous cycles. This allowed for the efficient follow-up of non-contact cases. 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.

Interviewers were instructed to make all reasonable attempts to obtain NPHS interviews with longitudinal respondents. For cases in which the timing of the interviewer's call (or visit) was inconvenient, an appointment was made to call back (or come back) at a more convenient time. If no one was home, numerous call backs were made. For individuals who refused to participate in the NPHS, a letter was sent from the Regional Office to the respondent, stressing the importance of the survey and the respondent's co-operation. This was followed by a second call (or visit) from the interviewer.

Refusals were followed up by senior interviewers, project supervisors or by other interviewers to try to convince respondents to participate in the survey. To maximise the response rate, a large number of non-response cases were also followed up in subsequent collection periods.

The failure to trace a longitudinal respondent was another type of non-response. Interviewers used several methods 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 other contacts, if collected in a previous cycle. In addition, interviewers were trained to follow up available public leads such as local telephone directories and directory assistance. If these leads were unsuccessful, the case was transmitted to an interviewer specially trained in tracing respondents. Tracer interviewers had access to Canada-wide telephone directories and reverse directories. The cumulative non-response rate due to failure to trace the longitudinal respondent is 5.4% of the total panel, which is relatively low for the seventh cycle of the survey. Section 9.2.3 presents non-response rates due to non-tracing with more details.

Attempts were made to contact panel members who moved within Canada or to the United States. For panel members living outside Canada and the United States, attempts were made to confirm their place of residence. The survey was not conducted if these members were still living outside Canada and outside the United States; however information was updated for the next cycle.

7. Data Processing

7.1 Editing

Editing was first performed on-line in the Computer Assisted Interview (CAI) application during data collection. It was not possible to enter out-of-range values, and flow errors were controlled by the skip pattern programmed in the CAI system. For example, CAI ensured that questions that did not apply to the respondent were not asked. In the case of contradictory responses between questions, warning messages were invoked. In some situations, the conflict had to be resolved before the interview could continue. In other situations, the contradiction was accepted and no corrective action had to be taken. Because of such cases, edits were developed to be performed after data collection at Head Office. Inconsistencies were usually corrected by setting one or more variables to "not stated". No imputation was performed.

7.2 Coding

For a number of questions in NPHS a long answer was captured (e.g. medication, illness, occupation, place of work). During data processing, the responses to these questions are coded according to the standard classification systems used by Statistics Canada.

Since the release of Cycle 4 microdata file, the industry and occupation data for all cycles, are coded to the North American Industrial Classification System (NAICS) and Standard Occupational Classification 1991 (SOC 1991) .

Since the release of Cycle 5 microdata file, the drug coding for all cycles is based on the Anatomical Therapeutic Chemical (ATC) classification developed by the World Health Organisation (WHO) as available on the Health Canada Drug Product Database (DPD) (September 2005). A complete list of the codes is available upon request.

Starting the release of Cycle 6 microdata file, the conditions or health problems causing activity restrictions for all cycles were coded based on the International Statistical Classification of Diseases and Related Health Problems, 10th Revision (ICD-10). Furthermore, the Musculoskeletal Impairment Supplementary Coding Scheme is not used anymore and it is not possible to code using ICD-9 at Statistics Canada. Therefore, the Cycle 7 microdata file and the following will only contain ICD-10 codes from cycle 1 to 7.

Many NPHS questions offers response categories and also allow to capture a long answer (specify). During data processing, and when appropriate, these long answers are coded in existing categories. Occasionally (mainly in Cycle 1) new categories were created.

The death of a longitudinal panel member is confirmed against the Canadian Vital Statistics Database – Deaths when possible. When the death is confirmed, the cause of death is capture and also the date of death (if different from the one collected during the survey). The cause of death is then coded using the ICD-10 only (since ICD-9 is not available anymore at Statistics Canada). Variables for panel members who died are set to "9" (i.e., not stated) in the dataset.

There is a total of 2,032 panel members who died during the first 7 cycles of NPHS. Among this number, 1,773 deaths have been confirmed with Canadian Vital Statistics Database –

Deaths. The year of death is known for 254 of the 259 remaining cases. Therefore, there are 5 cases left for which neither the year nor the cause of death is known. The following table presents the number of deaths and the number of confirmed death per year of death.

Table 7.A: Distribution of Deaths by Year of Death

Year of death	Number of deaths	Number of deaths confirmed with the Canadian Vital Statistics Database – Deaths
1994	25	25
1995	115	115
1996	152	149
1997	147	144
1998	188	184
1999	160	157
2000	180	176
2001	155	151
2002	182	181
2003	174	172
2004	191	188
2005	179	131
2006	146	0
2007	33	0
Unknown	5	0
TOTAL	2,032	1,773

7.3 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. For example, several variables may be combined to create a new derived variable. Derived variable names generally have a "D" in the fifth character of the variable name (see Section 11.3 for more detail on the variable naming convention). In other cases, see the document called "National Population Health Survey – Derived Variables Documentation – Cycles 1 to 7" for the details on how these variables were derived.

Grouped variables were created from certain variables; i.e. the values of the variable have been grouped in order to create another variable. In some cases, the derived variables are straightforward, involving collapsing response categories. Grouped variable names generally have a "G" in the fifth character of the variable name (see Section 11.3 for more detail on the variable naming conventions).

7.4 Estimation and 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 must be used to derive meaningful estimates from the survey. For example, if the number of individuals whose general health has deteriorated between two survey cycles 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.

The NPHS weighting method is presented in Chapter 8.

7.5 Definition of the Longitudinal Response Pattern (LONGPAT)

In each cycle, depending on the outcome of the interview, each member of the panel is assigned one of the following five statuses:

Completed (1): status given to panel members who provided a complete response to the interview (i.e., answered all the questions up to a given point in the NPHS questionnaire).

Deceased (2): status given to deceased panel members

Institutionalized (3): status given to members of the panel residing in health care institutions. This status indicates that a complete or partial response was obtained via the collection mechanism for the institutional component.

Partial response (4): status given to panel members who only partially responded to the NPHS questionnaire.

Non-response (5): status given to if none of the above status was assigned.

Over the cycles, the response statuses of a given respondent are concatenated into a single variable called the "longitudinal response pattern" (LONGPAT). This variable is available on the NPHS microdatafile and can be used to obtain rapidly the response profile of a member of the panel. This variable is also used to identify different analytical subsets as described in Section 7.7.

During data processing of a current cycle, an error is sometimes discovered in the response status of a previous cycle. Corrections are then made to the longitudinal response pattern. For example, a panel member with a response status of "non-response" the previous year is found to be deceased after having linked the data to the Canadian Vital Statistics Death Database. This person's response status is then set to "deceased" according to the date of death found in the database. For example, the variable LONGPAT for this person in cycle 6 was 115555 but it becomes 1155222 in cycle 7.

7.6 Definition of Full/Complete Response and Non-Response

Since Cycle 4, NPHS is strictly longitudinal. The definition of a response is not the same for longitudinal and cross-sectional purposes. For the NPHS longitudinal panel, a Full/Complete response includes panel members with the following statuses: completed, deceased and institutionalized.

Therefore, non-response includes panel members with the following statuses: partially completed and non-response.

7.7 Subsets of Respondents

In order to provide greater flexibility to users, a single microdata master file has been created for NPHS Cycle 7. This file includes all 17,276 NPHS panel members, regardless of their response patterns from Cycles 1 to 7. Within the master file, four subsets of respondents have been created along with corresponding sampling weights and the flags to make their identification easier. Refer to Chapter 8 for more information regarding the calculation of each subset's sampling weights and to Section 11.1 for the use of longitudinal weights. Table 7.B provides a description of the four subsets of respondents based on the type of response.

Table 7.B: Subsets of Respondents

Subset of Respondents	Type of Response	Flag	Number of Respondents
Longitudinal Square	Complete panel: all panel members regardless of their response pattern in Cycles 1 <u>to</u> 7.	None, all records	17,276
Longitudinal Full	All panel members with a complete response (Full) in Cycles 1 <u>to</u> 7.	WF6BLF	10,992
Longitudinal Square Share	All panel members regardless of their response pattern and who agreed to share their data.	WF6BSLS	16,007
Longitudinal Full Share	All panel members with a complete response (Full) in Cycles 1 <u>to</u> 7 and who agreed to share their data.	WF6BSLF	10,668

Users of the share file, provincial health departments, Health Canada and the Public Health Agency of Canada⁴, should note that the “Longitudinal Square Share” subset of respondents which includes the flag for the “Longitudinal Full Share” subset is provided separately on a CD-ROM with the corresponding sampling weights, for both subsets. The sampling weights and the flags of the other subsets are not on the share file CD-ROM.

⁴ The federal government developed a new structure to adapt to new needs in health. Health Canada and the new Public Health Agency of Canada (PHAC) come under the Minister of Health. The PHAC is one of the NPHS sharing partners.

8. Weighting

This chapter describes the weighting procedures for each subset of respondents described in Section 7.7. The longitudinal weighting process is necessarily different from that of cross-sectional weighting, for several reasons. First, longitudinal weights must represent the probability of selection of the unit of analysis at the time of sample selection. Since the longitudinal sample was selected in 1994/1995, the weights must reflect the probability of selecting the individual in Cycle 1 and not in subsequent cycles. In addition, the definition of a longitudinal response is different from that of a cross-sectional response, necessitating different adjustments particular to each type of non-response. Analysts should always use the longitudinal weights made from the subsets of respondents. The longitudinal weights have been calculated specifically to represent the 1994/1995 target population. In Cycles 1, 2 and 3, both cross-sectional and longitudinal files were produced. Although panel members were part of the cross-sectional and longitudinal files, their weights were not identical for these two types of files but rather adjusted to correctly represent the target population.

For Cycle 7, four sets of weights, WT64LS, WT6BLF, WT6BSLS and WT6BSLF have been created. Table 8.A shows the subsets of respondents and the corresponding sampling weights and flags. A panel member is part of a given subset when the flag is equal to 1.

Table 8.A: Subsets of Respondents and Corresponding Sampling Weights and Flags

Subset of respondents	Sampling Weight	Flag
Longitudinal Square	WT64LS	None, all records
Longitudinal Full	WT6BLF	WF6BLF
Longitudinal Square Share	WT6BSLS	WF6BSLS
Longitudinal Full Share	WT6BSLF	WF6BSLF

Only the WT6BLF, WT6BLFE, WT6BSLS and WT6BSLF weights have been adjusted for non-response. However, all four weights were post stratified to the 1994/1995 population estimates based on the 1996 Census counts by age group⁵ and sex within each province. Post-stratification is used to ensure that the four subsets of respondents represent correctly the 1994/1995 NPHS target population. The next section describes the NPHS longitudinal weighting method.

8.1 Longitudinal Weighting

The longitudinal weighting procedure is based on the weighting done for the Cycle 1 NPHS cross-sectional sample. Some weight adjustments were applied to the Cycle 1 cross-sectional weights in order to incorporate the additional sample used exclusively for cross-sectional purposes. These adjustments were removed for the longitudinal panel weight to create a “stripped” weight. This stripped weight is the starting point to obtain the longitudinal weight.

⁵ Post-stratification is done by using the updated date of birth instead of using the age variable at cycle 1 (DHC4_AGE) which is never updated.

8.1.1 Starting Point: Cycle 1 (1994/1995) Stripped Weights

The Cycle 1 stripped weights were obtained using the LFS basic weights as a starting point for all provinces except Quebec, where the basic weights from the “Enquête Sociale et de Santé” were taken as a starting point. Several adjustments were made to these weights to take into account the nature of the NPHS and to accurately represent the true probability of selection for each panel member. All of the adjustments that were made in Cycle 1 are kept for the subsequent cycles since the longitudinal sample always refers to the same population that is the population of 1994/1995.

A full description of the Cycle 1 weighting procedures still relevant for subsequent cycles is included in sections 11.3 and 11.4 of the Cycle 2 PUMF documentation. http://www.statcan.ca/english/sdds/document/3236_D7_T1_V2_E.pdf

From this point, adjustments were made to the stripped weight to obtain the various sets of longitudinal weights.

8.1.2 Adjustments to Create the Different Longitudinal Weights

8.1.2.1 Longitudinal Square Weight (WT64LS)

The longitudinal square weight **WT64LS** is to be used with the longitudinal square subset. It is calculated by post-stratifying the Cycle 1 stripped weight to the 1994/1995 population estimates based on 1996 Census counts by age group (0-11, 12-24, 25-44, 45-64, 65 and older) and sex within each province. The post-stratification adjustment is given by:

$$\frac{\text{Population estimate in a province/age/sex category}}{\text{Sum of "stripped" weights of respondent household members in a province/age/sex category}}$$

8.1.2.2 Longitudinal Full Weight (WT6BLF)

The longitudinal full subset includes only panel members with a full response, i.e. members who have a status of “complete”, “deceased” or “institutionalized” at each cycle. Panel members who are excluded from this subset were therefore non-respondents, i.e. they had a status of “partial response” or “non-response” at some point during the first seven cycles of the survey, and their weight must be redistributed to compensate for this non-response.⁶

The Cycle 1 stripped weight is the starting point and adjustments for non-response are made. A different non-response adjustment is made for each cycle, and these adjustments are cumulative from one cycle to another. For example, to obtain the Cycle 7 weights, the non-response adjustments for Cycles 2 to 7 are applied successively to the Cycle 1 stripped weights.

⁶ See section 7.5 for the definition of longitudinal response pattern and section 7.6 for the definition of full/complete for NPHS.

The adjustments necessary in order to obtain the Cycle 7 Longitudinal Full Weight are described below.

Adjustment 1: Adjustment for Cycle 2 (1996/1997) Non-Response

Adjusting for non-response was done using the weighting class approach. Weighting classes consist of groupings of respondents who 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 present units 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 Chi-Square Automatic Interaction Detection (CHAID) algorithm was used to determine the weighting classes. In order to produce more stable adjustments, a minimum of 30 units per weighting class was used.

Separate weighting classes were created for each province. Note that the province here refers to the province of residence at the time of the sample selection in 1994/1995. The Cycle 1 characteristics of the household as well as personal characteristics of the longitudinal member were considered. Some characteristics related to the sampling design of the survey or to the sampling weight were also considered in an effort to incorporate the sampling design of the survey into the analysis. Personal characteristics from the Health component were not used because they were not available for many longitudinal members in 1994/1995.

The variables chosen by the CHAID algorithm to build weighting classes to adjust for Cycle 2 non-response are listed in Table 8.B. Two variables from the Cycle 1 sample design were used: a flag indicating the presence of members under 25 years old in the household, and a flag indicating the presence of members under 12 years old in the household. The Cycle 1 non-response flag for income and the flag indicating if the individual was under age 16 were also used. Please refer to the Data Dictionary for a complete description of the variables listed in Table 8.B.

Table 8.B: Variables for Cycle 2 Non-Response Adjustment

DHC4_AGE	DHC4_MAR	GE34DURB	LFC4_1	SDC4DRAC
DHC4DECF	DHC4_OWEN	HCC4DMDC	RAC4F1	SDC4GCB7
DHC4_DWE	GE34DCMA	INC4DIA5	SDC4DAIM	SEX

To adjust for longitudinal members who did not respond in Cycle 2, the following adjustment is applied to the weight of respondents:

$$\frac{\textit{Sum of weights for all longitudinal members}}{\textit{Sum of weights for Cycles 1 and 2 responding longitudinal members}}$$

This adjustment is performed within each weighting class.

Adjustment 2: Adjustment for Cycle 3 (1998/1999) Non-Response

The 15,672⁷ records with a full response after two cycles, in other words those with a longitudinal response pattern 11, 12 or 13⁸, are taken as the starting point. A non-response adjustment is applied to the cases that have a status of “full response” after three cycles (defined as one of the following response patterns: 111, 112, 113, 122, 131, 132 or 133). All other response patterns, i.e. 114, 115, 134 and 135, are considered as non-responses. Records for which the panel member was deceased in Cycle 2 (pattern 122) or institutionalized since Cycle 2 (pattern 133) are treated differently from the rest. For these records, no non-response adjustment is made since their weight in Cycle 2 has been already adjusted to reflect the fact that some of the Cycle 2 non-respondents may have in fact been deceased or institutionalized.

Adjusting for non-response was done using the weighting class approach. Separate weighting classes were created for each province, i.e. the 1994/1995 province of residence. When adjusting for non-response in Cycle 3, only the Cycle 2 characteristics of the household as well as personal characteristics of the longitudinal member were considered. Again, as for Cycle 2, characteristics related to the sampling design of the survey or to the sampling weight were considered in an effort to incorporate the sampling design of the survey into the analysis. However, unlike for the Cycle 2 non-response adjustment, personal characteristics from the Health component were used, because they were available for all records that went into the Cycle 3 non-response adjustment.

The variables chosen by the CHAID algorithm to build weighting classes to adjust for Cycle 2 non-response are listed in table 8.C. A Cycle 1 sample design variable indicating the household type (adult or child) was also used, as well as a Cycle 2 item non-response flag for income. Please refer to the Data Dictionary for a complete description of the variables listed in Table 8.C.

⁷ When the Cycle 2 data were released, there were 15,670 records in the longitudinal full subset. With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now became so, and some others that were full/complete, were not anymore. Following these modifications, there were 15,672 cases with a full/complete response after two cycles instead of 15,670 cases.

⁸ See Section 7.6 for the definition of full/complete response for NPHS and Section 7.7 for the definition of longitudinal response pattern.

Table 8.C: Variables for Cycle 3 Non-Response Adjustment

AD_6_1	DHC6_AGE	INS6_4	SDC6DAIM	SMS6_9A
AD_6_7	DHC6_MAR	INS6_6	SDC6_4P	SMS6_13A
ALC6WKY	DV_6_65J	LFC6_41	SDC6_5A	SMS6_13C
ALC6_3	EDC6_3	MHC6DWK	SDC6_5F	SMS6_13E
AM56_SHA	ES_6_80	MHC6_1A	SDC6_6B	SMS6_16D
AM66_PXY	GE36LMOV	MHC6_1B	SDC6_7A	SMS6_18A
AM66_SHA	HCC6F1	MHC6_1F	SDC6_7B	SMS6_18D
BPC6_10	HSC6DPAD	MHC6_1L	SDC6_7D	SP36_CPA
CCC6DNUM	HWS6_5	MHC6_13	SEX	SSC6D2
CCC6_1L	INC6DIA5	PC_6_40	SHS6_4	SSC6_3
CCC6_1N	INC6_1A	RPC6_3	SMC6_2	SSS6_2
DGC6_1D	INC6_3B	RSS6_1	SMC6_5	SSS6_4

To adjust for longitudinal members who did not respond in Cycle 3, the following adjustment is applied to the weight of respondents:

$$\frac{\text{Sum of weights for Cycles 1 and 2 responding longitudinal members}}{\text{Sum of weights for Cycles 1, 2 and 3 responding longitudinal members}}$$

This adjustment is performed within each weighting class, and is calculated from records with the following longitudinal response patterns: 111 to 115, 131, 132, 134 and 135. Again, records for which the panel member was deceased in Cycle 2 or institutionalized since Cycle 2 are not part of this adjustment.

Adjustment 3: Adjustment for Cycle 4 (2000/2001) Non-Response

The 14,631⁹ records with a full response after three cycles are taken as the starting point. Once again, records for which the panel member was deceased in Cycle 2 or 3 or institutionalized since Cycle 2 or 3 are treated differently from the rest. For these records, no non-response adjustment is made since their weight in Cycle 2 or 3 has been already adjusted to reflect the fact that some of the Cycle 2 or Cycle 3 non-respondents may have in fact been deceased or institutionalized.

Here again, adjusting for non-response was done using the weighting class approach. Separate weighting classes were created for each design province i.e. the 1994/1995 province of residence. When adjusting for non-response in Cycle 4, only the Cycle 3 characteristics of the household as well as personal characteristics of the longitudinal member were considered. As for Cycle 3, characteristics related to the sampling design of the survey or to the sampling weight were considered in an effort to incorporate the sampling design of the

⁹ When the Cycle 3 data were released, there were 14,619 records in the longitudinal full subset. With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now became so, and some others that were full/complete, were not anymore. Following these modifications, there were 14,631 cases with a full/complete response after three cycles instead of 14,619 cases.

survey into the analysis. Personal characteristics from the Health component were used, because they were available for all records that went into the Cycle 4 non-response adjustment.

The variables chosen by the CHAID algorithm to build the weighting classes to adjust for Cycle 4 non-response are in Table 8.D. A Cycle 3 item non-response flag for income was also used. Please refer to the Data Dictionary for a complete description of the variables listed in table 8.D.

Table 8.D: Variables for Cycle 4 Non-Response Adjustment

CCC8DANY	DGC8_1A	HCC8_1	PAC8_1A	SDC8_6A
CCC8_1C	DHC8_AGE	INC8DIA5	PAC8_1J	SDC8_7A
CCC8_1L	DHC8DECF	ISC8_1	PY_8DH1	SEX
CCC8_1N	DHC8_OWN	NU_8_1B	RAC8F1	SSC8DEMO
CCC8_1V	FIC8F1	PAC8DFD	RPC8_2	SSC8DSOC
DGC8F1	GE38DURB	PAC8DLEI	SDC8_4A	TWC8_5

To adjust for longitudinal members who did not respond in Cycle 4, the following adjustment is applied to the weight of respondents:

$$\frac{\text{Sum of weights for Cycles 1, 2 and 3 responding longitudinal members}}{\text{Sum of weights for Cycles 1 to 4 responding longitudinal members}}$$

This adjustment is performed within each weighting class. Records for which the panel member was deceased in Cycle 2 or 3 or institutionalized since Cycle 2 or 3 are not part of this adjustment.

Adjustment 4: Adjustment for Cycle 5 (2002/2003) Non-response

The 13,597¹⁰ records with a full response after four cycles are taken as the starting point. Once again, records for which the panel member was deceased in Cycle 2, 3 or 4 or institutionalized since Cycle 2, 3 or 4 are treated differently from the rest. For these records, no non-response adjustment is made since their weight in Cycle 2, 3 or 4 has been already adjusted to reflect the fact that some of the Cycle 2, Cycle 3 or Cycle 4 non-respondents may have in fact been deceased or institutionalized.

Here again, adjusting for non-response was done using the weighting class approach. Separate weighting classes were created for each design province i.e. the 1994/1995 province of residence. When adjusting for non-response in Cycle 5, only the Cycle 4 characteristics of the household as well as personal characteristics of the longitudinal member were considered. As for Cycle 4,

¹⁰ When the Cycle 4 data were released, there were 13,582 records in the longitudinal full subset. With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now, became so, and some others that were full/complete, were not anymore. Following these modifications, there were 13,597 cases with a full/complete response after four cycles instead of 13,582 cases.

characteristics related to the sampling design of the survey or to the sampling weight were considered in an effort to incorporate the sampling design of the survey into the analysis. Personal characteristics from the Health component were used, because they were available for all records that went into the Cycle 5 non-response adjustment.

The variables chosen by the CHAID algorithm to build the weighting classes to adjust for Cycle 5 non-response are in Table 8.E. Three Cycle 1 design variables were also used, one indicating the presence of household members under the age of 12, the other indicating the presence of household members under the age of 25, and the last one indicating the household type (adult or child). A Cycle 4 item non-response flag for income was also used. Please refer to the Data Dictionary for a complete description of the variables listed in table 8.E.

Table 8.E: Variables for Cycle 5 Non-Response Adjustment

ALC0_3	DHC0_OWN	IMM	MHC0_1J	SMC0_2
ALC0DTYP	DHC0DL12	INC0DIA5	MHC0DCH	ST_0DC4
ALC0DWKY	DHC0DLE5	ISC0_1	MHC0DDS	ST_0DC5
AM60_SHA	GE30DURB	LSC0_1	PAC0DFD	ST_0DC6
BPC0_10	GHC0_21	LSC0DPFT	PAC0DLEI	ST_0DC8
CCC0DANY	HCC0DHPC	MHC0_16	SDC0_4A	ST_0DR2
DGC0F1	HSC0DHSI	MHC0_1A	SDC0_6A	ST_0DW3
DHC0_AGE	HWC0DSW	MHC0_1F	SEX	ST_0DW6

To adjust for longitudinal members who did not respond in Cycle 5, the following adjustment is applied to the weight of respondents:

$$\frac{\text{Sum of weights for Cycles 1 to 4 responding longitudinal members}}{\text{Sum of weights for Cycles 1 to 5 responding longitudinal members}}$$

This adjustment is performed at the weighting class level. Records for which the panel member was deceased in Cycle 2, 3 or 4 or institutionalized since Cycle 2, 3 or 4 are not part of this adjustment.

Adjustment 5: Adjustment for Cycle 6 (2004/2005) Non-Response

The 12,559¹¹ records with a full response after five cycles are taken as the starting point. Once again, records for which the panel member was deceased in Cycle 2, 3, 4 or 5 or institutionalized since Cycle 2, 3, 4 or 5 are treated differently from the rest. For these records, no non-response adjustment is made since their weight in Cycle 2, 3, 4 or 5 has been already adjusted to reflect the fact that some

¹¹ When the Cycle 5 data were released, there were 12,546 records in the longitudinal full subset. With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now, became so, and some others that were full/complete, were not anymore. Following these modifications, there were 12,559 cases with a full/complete response after five cycles instead of 12,546 cases.

of the Cycle 2, 3, 4 or 5 non-respondents may have in fact been deceased or institutionalized.

Here again, adjusting for non-response was done using the weighting class approach. Separate weighting classes were created for each design province i.e. the 1994/1995 province of residence. When adjusting for non-response in Cycle 6, only the Cycle 5 characteristics of the household as well as personal characteristics of the longitudinal member were considered. As for Cycle 5, characteristics related to the sampling design of the survey or to the sampling weight were considered in an effort to incorporate the sampling design of the survey into the analysis. Personal characteristics from the Health component were used, because they were available for all records that went into the Cycle 6 non-response adjustment.

The variables chosen by the CHAID algorithm to build the weighting classes to adjust for Cycle 6 non-response are in Table 8.F. Two variables from the Cycle 1 sample design were used: a flag indicating the presence of household members less than 12 years of age and a flag indicating the presence of household members less than 25 years of age. A Cycle 5 item non-response flag for income was also used. Please refer to the Data Dictionary for a complete description of the variables listed in table 8.F.

Table 8.F: Variables for Cycle 6 Non-Response Adjustment

ALC2_2	DHC2DL12	LSC2_1	PAC2_3E	SSC2DTNG
ALC2_3	EDC2_4	LSC2_21	PAC2_3F	ST_2DC5
ALC2DWKY	EDC2D2	LSC2DPFT	PAC2DFD	ST_2DC6
CCC2_1C	GE32DURB	MHC2_1A	PAC2DFM	ST_2DC7
CCC2_1F	HSC2DEMO	MHC2_1C	RAC2_1C	ST_2DC8
CCC2_1L	HWC2DISW	MHC2_1F	RAC2_6D	ST_2DC9
CCC2DANY	HCC2DMDC	MHC2_1G	SEX	ST_2DW1
CCC2DNUM	IMM	PAC2_1F	SMC2_2	ST_2DW2
DHC2_AGE	INC2_1A	PAC2_3A	SMC2DTYP	ST_2DW5
DHC2_OWN	INC2DHH	PAC2_3B	SMC2DYRS	ST_2DW6
DHC2DHSZ	INC2DIA5			

To adjust for longitudinal members who did not respond in Cycle 6, the following adjustment is applied to the weight of respondents:

$$\frac{\text{Sum of weights for Cycles 1 to 5 responding longitudinal members}}{\text{Sum of weights for Cycles 1 to 6 responding longitudinal members}}$$

This adjustment is performed at the weighting class level. Records for which the panel member was deceased in Cycle 2, 3, 4 or 5 or institutionalized since Cycle 2, 3, 4 or 5 are not part of this adjustment.

Adjustment 6: Adjustment for Cycle 7 (2006/2007) Non-Response

The 11,619¹² records with a full response after six cycles are taken as the starting point. Once again, records for which the panel member was deceased in Cycle 2, 3, 4, 5 or 6 or institutionalized since Cycle 2, 3, 4, 5 or 6 are treated differently from the rest. For these records, no non-response adjustment is made since their weight in Cycle 2, 3, 4, 5 or 6 has been already adjusted to reflect the fact that some of the Cycle 2, 3, 4, 5 or 6 non-respondents may have in fact been deceased or institutionalized.

Here again, adjusting for non-response was done using the weighting class approach. Separate weighting classes were created for each design province i.e. the 1994/1995 province of residence. When adjusting for non-response in Cycle 7, only the Cycle 6 characteristics of the household as well as personal characteristics of the longitudinal member were considered. As for Cycle 6, characteristics related to the sampling design of the survey or to the sampling weight were considered in an effort to incorporate the sampling design of the survey into the analysis. Personal characteristics from the Health component were used, because they were available for all records that went into the Cycle 7 non-response adjustment.

The variables chosen by the CHAID algorithm to build the weighting classes to adjust for Cycle 7 non-response are in Table 8 G. Two variables from the Cycle 1 sample design were used: a flag indicating the presence of household members less than 25 years of age and a flag indicating the household type (adult or child). A Cycle 6 age group variable was also used. Please refer to the Data Dictionary for a complete description of the variables listed in table 8.G.

Table 8.G: Variables for Cycle 7 Non-Response Adjustment

ALCA_2	DHCA_MAR	INCADHH	LSCA_21	SSCADAFF
ALCA_3	EDCA_4	INCADIAG5	PACADFD	STCADW1
DHCADECF	GHCA_2	INCA_3B	SDCA_6A	STCADW6
DHCADLVG	IMM	LSCADPFT	SMCA_2	STCADW7
DHCA_AGE				

To adjust for longitudinal members who did not respond in Cycle 7, the following adjustment is applied to the weight of respondents:

$$\frac{\text{Sum of weights for Cycles 1 to 6 responding longitudinal members}}{\text{Sum of weights for Cycles 1 to 7 responding longitudinal members}}$$

¹² When the Cycle 6 data were released, there were 11,593 records in the longitudinal full subset. With the information from Cycle7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now, became so, and some others that were full/complete, were not anymore. Following these modifications, there were 11,619 cases with a full/complete response after six cycles instead of 11,593 cases.

This adjustment is performed at the weighting class level. Records for which the panel member was deceased in Cycle 2, 3, 4, 5 or 6 or institutionalized since Cycle 2, 3, 4, 5 or 6 are not part of this adjustment.

Adjustment 7: Post-Stratification Adjustment

The weight of the units that are part of the full subset was post-stratified to the 1994/1995 population estimates based on 1996 Census counts by age group (0-11, 12-24, 25-44, 45-64, 65 and older) and sex within each province. This is done to ensure that the 1994/1995 population is accurately represented in any estimates produced from the longitudinal file. This adjustment is given by:

$$\frac{\text{Population estimate in a province/age/sex category}}{\text{Sum of weights of Cycles 1 to 7 responding longitudinal members in a province/age/sex category}}$$

The final longitudinal weight **WT6BLF** is calculated by taking the Cycle 1 stripped weight and multiplying that value by adjustments 1 to 7.

8.1.2.3 Longitudinal Square Share Weight (WT6BSLS)

Creation of the share square subset was started at Cycle 5. This subset includes the panel members who agreed to share the information provided from all interviews conducted as part of NPHS with provincial ministries of health, Health Canada and Public Health Agency of Canada. As these partners only receive the records of these sharers, a special weight, WT6BSLS, must be derived so that estimates computed from this subset correctly represent the total population.

To calculate this weight, the following two adjustments were applied to the Cycle 1 stripped weight: a non-response adjustment (do not agree to share) and a post-stratification adjustment.

Adjustment 1: Non-Response Adjustment (do not agree to share)

As for the full weight computation (section 8.1.2.2), adjusting for the non-response was done using the weighting class approach. Weighting classes consist of groupings of panel members who share the same propensity to agree to share their information. The characteristics of the household as well as personal characteristics of the longitudinal panel members from each of the seven cycles were considered to define membership in the weighting classes. The Chi-Square Automatic Interaction Detection (CHAID) algorithm was used to determine the weighting classes and a minimum of 30 units per weighting class was used in order to produce more stable adjustments. Separate weighting classes were created for each province, i.e. the 1994/1995 province of residence.

The variables chosen by the CHAID algorithm to build weighting classes for the non-response adjustment are listed in table 8.H. Variables identifying if the panel member had a status of “complete”, “deceased”, “institutionalized”, “partial response” or “non-response” to Cycles 3, 4, 6 and 7 were used as well as a flag indicating whether the panel member was part of the full subset. One variable

from the Cycle 1 sample design representing a flag indicating the presence of household members under 25 years of age was also used as well as flags indicating non-response to the income module in Cycles 3 and 7. Please refer to the Data Dictionary for a complete description of the variables listed in table 8.H.

Table 8.H: Variables for the Non-Response Adjustment

ALCA_2	DHC4DSZ	DHCB_OWN	EDCBD2	INCBDIAG5
DHC0DECF	DHC8_MAR	EDC2D2	IMM	SDC4_6A
DHC4_AGE	DHCB_MAR	EDC8D2	INCBDDH	SEX

To adjust for longitudinal members who refused to share their data in Cycle 7, the following adjustment is applied to the weight of the sharers:

$$\frac{\text{Sum of weights for all longitudinal members}}{\text{Sum of weights for longitudinal members agreeing to share}}$$

This adjustment is performed within each weighting class.

Adjustment 2: Post-Stratification Adjustment

The weight of the units that are part of the square share subset was post-stratified to the 1994/1995 population estimates based on 1996 Census counts by age group (0-11, 12-24, 25-44, 45-64, 65 and older) and sex within each province. This is done to ensure that the 1994/1995 population is accurately represented in any estimates produced from the longitudinal file. This adjustment is given by:

$$\frac{\text{Population estimate in a province/age/sex category}}{\text{Sum of weights for longitudinal members agreeing to share in Cycle 7 in a province/age/sex category}}$$

The final weight for the square share subset, **WT6BSLS**, is calculated by taking the Cycle 1 stripped weight and multiplying this value by adjustments 1 and 2.

8.1.2.4 Longitudinal Full Share Weight (WT6BSLF)

As for the square share subset, the full share subset includes the panel members who agreed to share the information provided from all interviews conducted as part of NPHS with provincial ministries of health, Health Canada and Public Health Agency of Canada but only those who had a full response in Cycles 1 to 7. As these partners only receive the records of these sharers, a special weight must be derived so that estimates computed from this subset correctly represent the total population.

A simple adjustment is made to the longitudinal full weight to create the full share weight. This adjustment is given by:

$$\frac{\text{Sum of weights for Cycles 1 to 7 responding longitudinal members in a province / longitudinal pattern / age-sex category}}$$

Sum of weights for Cycles 1 to 7 responding longitudinal members who agreed to share, in a province / longitudinal pattern / age-sex category

Note that in Cycles 3 to 7, a few of the original longitudinal response patterns were collapsed in order to produce more stable adjustments. The grouping was done for a few province/age-sex categories that had few observations in some of the longitudinal patterns representing deceased or institutionalized. In each case, the problematic response pattern was grouped with another longitudinal pattern in the same province/age-sex category, so that the sum of the weights would still give the correct population counts. The final longitudinal share weight, **WT6BSLF**, is obtained by multiplying the longitudinal full weight, WT6BLF, by this adjustment. Note that since this adjustment is done with respect to the post-stratification classes, no additional post-stratification is necessary.

9. Data Quality

Data quality is an important aspect for any survey. Examining data quality allows the verification of the reliability and accuracy of the data collected, as well as help to determine what steps should be taken to improve data quality in future cycles.

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 misunderstand the questions asked, the answers may be incorrectly entered or errors may be introduced in the processing and tabulation of the data. These are all examples of non-sampling errors.

9.1 Sampling Errors

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

For example, suppose hypothetically that one estimates that 25% of Canadians aged 12 and over have experienced an improvement in their general health between Cycle 1 and Cycle 2 of the survey and that this estimate is found to have a standard deviation of .003. Then the CV of the estimate is calculated as:

$$(.003/.25) \times 100\% = 1.20\%.$$

Statistics Canada commonly uses CV results to verify the quality of statistical estimates produced when analyzing data, and strongly urges users producing estimates from NPHS data files to also do so. For guidelines on how to interpret CV results, see the table at the end of Section 10.4.

9.1.1 Bootstrap Method for Variance Estimation

In order to determine the quality of the estimate and to calculate the CV, the standard deviation must be calculated. Confidence intervals as well as a large number of statistical tests also require the standard deviation of the estimate.

The NPHS uses a multi-stage survey design, which means that there is no simple formula that can be used to calculate variance estimates. Therefore, an approximate method was needed. The bootstrap method is used because the

sample design information needs to be taken into account when calculating variance estimates. The bootstrap method does this, and with the use of the Bootvar program, remains a method that is fairly easy for users to use.

The bootstrap re-sampling method used in the NPHS involves the selection of simple random samples known as replicates, and the calculation of the variation in the estimates from replicate to replicate. In each stratum, a simple random sample of (n-1) of the n clusters is selected with replacement to form a replicate. Note that since the selection is with replacement, a cluster may be chosen more than once. In each replicate, the survey weight for each record in the (n-1) selected clusters is recalculated. These weights are then post-stratified according to demographic information in the same way as the sampling design weights in order to obtain the final bootstrap weights.

The entire process (selecting simple random samples, recalculating and post-stratifying weights for each stratum) is repeated B times, where B is large. The NPHS typically uses B=500, to produce 500 bootstrap weights. To obtain the bootstrap variance estimator, the point estimate for each of the B samples must be calculated. The standard deviation of these estimates is the bootstrap variance estimator. Statistics Canada has developed a program that can perform all of these calculations for the user: the Bootvar program. For more information on Bootstrap weights, please refer to Section 11.2.

The Bootvar program is available in both SAS and SPSS formats. It is made up of macros that compute variances for totals, ratios, differences between ratios and for linear and logistic regression.

The Bootvar program is provided with bootstrap weights and a document explaining how to modify and use the program to suit user's needs.

9.2 Non-Sampling Errors

Considerable time and effort was made to reduce non-sampling errors in the NPHS. Quality assurance measures were implemented at each step of data collection and processing to monitor the quality of the data. These measures included the use of highly skilled interviewers, extensive training with respect to the survey procedures and questionnaire, and the observation of interviewers to detect problems and give solution if needed. Testing of the CAI application and field tests were also essential 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 is minimal; once the questionnaire is started, it tends to be completed with very little non-response. In most cases, partial non-response to the survey occurred when the respondent refused to answer a question, could not recall the requested information, or could not provide personal or proxy information. Total non-response occurred because it was impossible to trace or reach 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 persons who responded to the survey to

compensate for those who did not respond. See Section 8.1.2 for details of the weight adjustment for non-response.

This section presents some information dealing with different aspects of non-response. Discussed first is overall non-response, where response rates from each cycle are presented. This is followed by sections related to refusals, panel member untraced and attrition. Finally, item non-response is briefly examined.

9.2.1 Response Rates

This chapter presents the response rates and describes how they are computed. The calculation of Cycle 1 response rates is not the same as the calculation of the response rates for the other cycles. Cycle 1 response rates are based on the 20,095 in-scope persons selected to form the panel while response rates for subsequent cycles are based on the 17,276 individuals who form the longitudinal panel. Another important difference: for the first three cycles, the selected-person response rate is calculated both for the General component (H05) and for the Health component (H06) (see section 4.1). Since the survey became purely longitudinal in Cycle 4 and there was no longer a distinction between these two components, there is only one longitudinal panel response rate since cycle 4.

9.2.1.1 Cycle 1 (1994/1995) Response Rates

Cycle 1 response rates are based on the 20,095 in-scope persons selected to form the panel. Consequently, persons who were part of the 3,165 out-of-scope households (status code = 017, 018, 023, 024)¹³ and the 2,983 excluded households by the rejective method (see section 5.1.2) were excluded from the panel and from the calculations of the Cycle 1 response rates.

Selected-Person Response Rate for H05

$$\frac{\text{\# of selected persons responding to the H05 component}}{\text{all in-scope selected persons}}$$

The selected-person response rate for the H05 component at the Canada level for the NPHS was **86.0%**. At the provincial level, this rate varied from 80.7% in Ontario to 91.0% in Alberta.

Selected-Person Response Rate for H06

$$\frac{\text{\# of selected persons responding to the H06 component}}{\text{all in-scope selected persons}}$$

The selected-person response rate for the H06 component was **83.6%** at the Canada level, and ranged from 77.8% in Ontario to 89.1% in Alberta.

¹³ 017 = Other ineligible dwelling (e.g., embassy).
018 = Rejected household.
023 = Under construction or demolished.
024 = Vacant dwelling.

Relevant information for the calculation of response rates is given in Table 9.A and response rates at the national and provincial level are presented in Table 9.B.

Table 9.A: Relevant Information for Calculation of Response Rates for Cycle 1

Cycle 1 (1994/1995)				
Number of in-scope selected persons (1) = (2)+(4) or (3)+(5)	Number of respondents		Number of non-respondents	
	H05 (2)	H06 (3)	H05 (4)	H06 (5)
20,095	17,276	16,794	2,819	3,301

Table 9.B: Cycle 1 Response Rates

Province	Cycle 1 Response Rates (1994/1995)	
	H05 (2) / (1)	H06 (3) / (1)
Newfoundland/Labrador*	89.3%	86.9%
Prince Edward Island	87.6%	84.9%
Nova Scotia	85.4%	82.1%
New Brunswick	88.1%	86.0%
Quebec	87.5%	85.6%
Ontario	80.7%	77.8%
Manitoba	89.5%	87.0%
Saskatchewan	88.5%	86.9%
Alberta	91.0%	89.1%
British Columbia	85.2%	82.8%
Canada	86.0%	83.6%

* Labrador was not part of Newfoundland when sample was selected.

9.2.1.2 Cycle 2 (1996/1997) and Cycle 3 (1998/1999) Response Rates

All Cycle 2 and Cycle 3 response rates are based on the 17,276 individuals who form the longitudinal panel. Persons with a status “*deceased*” or “*institutionalized*” are counted as a response for longitudinal purposes. However, persons with a status “*partial response*” are counted as a non-response (see section 7.6). No panel members are classified as out-of-scope.

Panel response rate for H05

$$\frac{\text{\# of panel members responding to the H05 component or who have died or been institutionalized}}{\text{\# of panel members}}$$

At the Canada level, the panel response rate for the H05 component was **93.6%** in Cycle 2 and **88.9%** in Cycle 3. At the provincial level, this rate varied from 90.4% in British Columbia to 96.1% in Newfoundland in Cycle 2 and from 84.2% in British Columbia to 92.5% in Newfoundland in cycle 3.

Panel response rate for H06

$$\frac{\text{\# of panel members responding to the H06 component or who have died or been institutionalized}}{\text{\# of panel members}}$$

At the Canada level, the panel response rate for the H06 component was **92.8%** in Cycle 2 and **88.3%** in Cycle 3. At the provincial level, this rate varied from 89.8% in British Columbia to 94.8% in Newfoundland in Cycle 2 and from 83.9% in British Columbia to 92.0% in Newfoundland in cycle 3.

Relevant information for the calculation of response rates is given in Table 9.C and response rates at the national and provincial level are presented in Table 9.D.

Table 9.C: Relevant Information for Calculation of Response Rates for Cycles 2 and 3

Cycle	Number of Panel Members ¹⁴					
	Deceased (1)	Institutionalized (2)	Complete		Non-respondent	
			H05 (3)	H06 (4)	H05 (5)	H06 (6)
2	289	62	15,819	15,687	1,112	1,238
3	615	114	14,647	14,532	1,912	2,015

¹⁴ With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now became so, and some others that were full/complete, were no longer. Numbers presented in this table show the situation after the modifications.

Table 9.D: Panel Response Rate for Cycles 2 and 3

Province	Response rate			
	Cycle 2 (1996/1997)		Cycle 3 (1998/1999)	
	H05 <u>(1) + (2) + (3)</u> 17 276	H06 <u>(1) + (2) + (4)</u> 17 276	H05 <u>(1) + (2) + (3)</u> 17 276	H06 <u>(1) + (2) + (4)</u> 17 276
Newfoundland/Labrador*	96.1%	94.8%	92.5%	92.0%
Prince Edward Island	95.0%	94.3%	91.3%	90.7%
Nova Scotia	94.6%	94.0%	90.1%	89.4%
New Brunswick	94.8%	94.3%	89.5%	89.2%
Quebec	95.1%	94.1%	89.3%	88.3%
Ontario	92.1%	91.4%	87.6%	86.9%
Manitoba	95.4%	94.4%	91.1%	90.6%
Saskatchewan	94.7%	94.0%	91.0%	91.0%
Alberta	91.8%	91.5%	88.7%	88.4%
British Columbia	90.4%	89.8%	84.2%	83.9%
Canada	93.6%	92.8%	88.9%	88.3%

* Labrador was not part of Newfoundland when sample was selected.

9.2.1.3 Cycle 4 (2000/2001), Cycle 5 (2002/2003), Cycle 6 (2004/2005) and Cycle 7 (2006/2007) Response Rates

As for Cycles 2 and 3, the Cycles 4, 5, 6 and 7 longitudinal response rates are based on the 17,276 members of the longitudinal panel. Persons with a status “*deceased*” or “*institutionalized*” are counted as a response for longitudinal purposes. However, persons with a status “*partial response*” are counted as a non-response (see section 7.6). No panel members are classified as out-of-scope. As of Cycle 4, NPHS is now purely longitudinal and no longer distinguishes the H05 questionnaire from the H06 questionnaire; only one response rate is calculated and is equivalent to the H06.

Response rate

$$\frac{\text{\# of panel members responding or who have died or been institutionalized}}{\text{\# of longitudinal panel members}}$$

At the Canada level, the panel member response rate was **84.9%** in Cycle 4, **80.8%** in Cycle 5 and **77.6%** in Cycle 6. In Cycle 7, the response rate was **77.0%**.

Relevant information for the calculation of response rates is given in Table 9.E and response rates at the national and provincial level are presented in Table 9.F.

Table 9.E: Relevant Information for Calculation of Response Rates for Cycles 4 to 7

Cycle	Number of Panel Members ¹⁵			
	Deceased (1)	Institutionalized (2)	Complete (3)	Non- respondent (4)
4	977	133	13,560	2,606
5	1,311	161	12,483	3,321
6	1,680	144	11,590	3,862
7	2,032	148	11,119	3,977

Table 9.F: Panel Response Rate for Cycles 4, 5, 6 and 7

Province	Response Rate = $\frac{(1) + (2) + (3)}{17276}$			
	Cycle 4	Cycle 5	Cycle 6	Cycle 7
Newfoundland/Labrador*	89.2%	82.6%	80.9%	82.6%
Prince Edward Island	88.2%	84.1%	80.6%	82.5%
Nova Scotia	87.0%	81.6%	81.0%	82.4%
New Brunswick	83.8%	79.4%	75.8%	78.3%
Quebec	85.5%	79.8%	75.2%	77.3%
Ontario	82.3%	80.1%	75.0%	72.6%
Manitoba	89.0%	82.7%	81.6%	78.4%
Saskatchewan	90.9%	84.8%	82.9%	81.0%
Alberta	83.3%	80.3%	79.3%	77.6%
British Columbia	80.8%	77.9%	76.5%	73.5%
Canada	84.9%	80.8%	77.6%	77.0%

* Labrador was not part of Newfoundland when sample was selected.

¹⁵ With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now became so, and some others that were full/complete, were no longer. Numbers presented in this table show the situation after the modifications.

9.2.2 Refusal Rates

In a longitudinal survey, non-response to one or more cycles is quite costly. In a sense, it breaks the chronological sequence of information available on a respondent, which may make it more complex to conduct analyses on the data. Worse yet, chronic non-response to the survey reduces the sample size available for analysis, thus diminishing the potential for observing statistically significant results. Consequently, many efforts have been made in the NPHS to minimize non-response.

Despite all efforts made to convert refusals (see section 6.4), they remain the most substantial source of non-response for the NPHS. Even though the intention is to follow all 17,276 panel members over time, not all records (panel members) are sent out for collection each cycle, such as the categorical refusals. Note that when the panel member has been confirmed dead through a match to the mortality files or the year of death of the panel member is known, the case is considered complete for the remaining duration of the survey, and is no longer sent out for collection.

Two different refusal rates for each cycle can be calculated, one based only on those panel members that were sent out for collection, and the other based on all 17,276 panel members. It can be seen in Table 9.G, which displays both of these rates for Cycles 2 to 7, that both refusal rates increased with each cycle but seem to stabilize over time. However, it must be noted that efforts to convert refusals are well worthwhile because some panel members that had refused for a few consecutive cycles, have ended up participating in subsequent cycles of the survey.

Table 9.G: Refusal Rates by Cycle¹⁶

Cycle	Number of panel members that went out	Number of refusal during collection	Refusal rate based on panel members sent out	Number of refusals that were not sent out	Total number of refusals	Refusal rate based on all 17,276 panel members
2	17,266	538	3.1%	1	539	3.1%
3	16,582	601	3.6%	460	1,061	6.2%
4	16,186	1,014	6.3%	514	1,528	8.9%
5	15,616	1,301	8.3%	664	1,965	11.5%
6	14,743	1,071	7.3%	1,213	2,284	13.4%
7	13,857	844	6.1%	1,444	2,288	13.2%

¹⁶ With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Numbers presented in this table show the situation after the modifications.

9.2.3 Unable to Trace Rates

After refusal, the failure to trace a longitudinal panel member is the second most substantial source of non-response for the NPHS. Despite the numerous efforts from the interviewers (discussed in Section 6.4), the cumulative unable-to-trace rate is increasing with the passing cycles but many attempts were put in place to keep this rate as low as possible. Two different unable to trace rates for each cycle can be calculated, one based only on those panel members that were sent out for collection, and the other based on all 17,276 panel members. It can be seen in Table 9.H, which displays both of these rates for Cycles 2 to 7 that both unable to trace rates increased with each cycle.

Table 9.H: Unable to Trace Rates by Cycle¹⁷

Cycle	Number of panel members that went out	Number of unable to trace during collection	Unable to trace rate based on panel members sent out	Number of unable to trace that were not sent out	Total number of unable to trace	Unable to trace rate based on all 17,276 panel members
2	17,266	295	1.7%	0	295	1.7%
3	16,582	359	2.2%	0	359	2.1%
4	16,186	500	3.1%	0	500	2.9%
5	15,616	698	4.5%	0	698	4.0%
6	14,743	869	5.9%	2	871	5.0%
7	13,857	686	5.0%	247	933	5.4%

9.2.4 Attrition Rates

In a longitudinal survey, attrition is a loss in sample size due to non-respondents. For the first five cycles, the attrition of the NPHS sample was defined only by whether or not a panel member was part of the full subset. Therefore, when a non-response was observed for a panel member, it was considered part of attrition. Defining attrition in this way permitted users to better understand this subset of data but it painted a somewhat pessimistic portrait of the sample's actual state. Attrition can be defined in different ways depending on the subset of data and the method of analysis used. This section presents the attrition based on the Full subset of data. Another way to describe attrition is also discussed in section 13 where the method called cycle twinning is explained. This method reduces the impact of attrition on the data. For more information on the different subsets of data, please see section 7.7.

¹⁷ With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Numbers presented in this table show the situation after the modifications.

9.2.4.1 Attrition Rates Based on the Longitudinal Full Subset

Here, attrition has been defined by whether or not a panel member is part of the full subset. Therefore, when a non-response is observed for a panel member, it is considered part of attrition. Two different attrition rates are calculated: one showing the attrition rate between two consecutive cycles, the other showing the cumulative attrition rate based on the original sample. Both of these rates are calculated using the number of individuals found in the Full subset of respondents.

Relevant information for calculation of attrition rates:

Number of longitudinal panel members:	17,276
Number of individuals in the Cycle 2 Full subset:	15,672 ¹⁸
Number of individuals in the Cycle 3 Full subset:	14,631 ¹⁹
Number of individuals in the Cycle 4 Full subset:	13,597 ²⁰
Number of individuals in the Cycle 5 Full subset:	12,559 ²¹
Number of individuals in the Cycle 6 Full subset:	11,619 ²²
Number of individuals in the Cycle 7 Full subset:	10,992

Attrition rates between two cycles:

Cycle 2 (1996/1997):	$\frac{17,276 - 15,672}{17,276} = \frac{1,604}{17,276} = 9.3\%$
Cycle 3 (1998/1999) :	$\frac{15,672 - 14,631}{15,672} = \frac{1,041}{15,672} = 6.6\%$
Cycle 4 (2000/2001) :	$\frac{14,631 - 13,597}{14,631} = \frac{1,034}{14,631} = 7.1\%$
Cycle 5 (2002/2003) :	$\frac{13,597 - 12,559}{13,597} = \frac{1,038}{13,597} = 7.6\%$
Cycle 6 (2004/2005) :	$\frac{12,559 - 11,619}{12,559} = \frac{940}{12,559} = 7.5\%$
Cycle 7 (2006/2005) :	$\frac{11,619 - 10,992}{11,619} = \frac{627}{11,619} = 5.4\%$

¹⁸ When the Cycle 2 data were released, there were 15,670 records in the longitudinal full subset. With the information from Cycle 7, we went back to previous cycles and could confirm, among other things, that some non-responses were in fact a death and some deaths were a non-response. Some cases that were not full/complete until now became so, and some others that were full/complete, were no longer. Following these modifications, there are 15,672 cases with a full response after two cycles instead of 15,670 cases.

¹⁹ For the same reasons as the previous note, after the modifications, there are 14,631 records in the Cycle 3 longitudinal full subset instead of 14,619.

²⁰ For the same reasons as the previous note, after the modifications, there are 13,597 records in the Cycle 4 longitudinal full subset instead of 13,582

²¹ For the same reasons as the previous note, after the modifications, there are 12,559 records in the Cycle 5 longitudinal full subset instead of 12,546.

²² For the same reasons as the previous note, after the modifications, there are 11,619 records in the Cycle 6 longitudinal full subset instead of 11,593.

Cumulative Attrition rates for the longitudinal full subset:

Cycle 2 (1996/1997) :	$\frac{17,276-15,672}{17,276}$	=	$\frac{1,604}{7,276}$	=	9.3%
Cycle 3 (1998/1999) :	$\frac{17,276-14,631}{17,276}$	=	$\frac{2,645}{7,276}$	=	15.3%
Cycle 4 (2000/2001) :	$\frac{17,276-13,597}{17,276}$	=	$\frac{3,679}{17,276}$	=	21.3%
Cycle 5 (2002/2003) :	$\frac{17,276-12,559}{17,276}$	=	$\frac{4,717}{17,276}$	=	27.3%
Cycle 6 (2004/2005) :	$\frac{17,276-11,619}{17,276}$	=	$\frac{5,657}{17,276}$	=	32.7%
Cycle 7 (2006/2007) :	$\frac{17,276-10,992}{17,276}$	=	$\frac{6,284}{17,276}$	=	36.4%

As is typically the case in longitudinal surveys, the attrition rate between Cycles 1 and 2 is considerably higher (9.3%) than those subsequently observed. The subsequent attrition rates are more constant between cycles. Cumulatively, after seven cycles, about 36% of the panel has eroded based on the full subset. Table 9.I presents the most important attrition type by cycle. As already mentioned, refusal and unable to trace are still the most substantial sources of attrition.

Table 9.I: Attrition Type by Cycle – Longitudinal Full Subset of Respondents

Attrition type	Attrition cycle by cycle						Cumulative Attrition
	Cycle 2	Cycle 3	Cycle 4	Cycle 5	Cycle 6	Cycle 7	Cycle 7
Refusal	30,6%	45,4%	52,2%	51,6%	48,0%	41,3%	43,8%
Unable to trace	17,1%	12,5%	16,8%	15,6%	20,3%	21,7%	17,0%
Partial response in Cycle 1	30,0%	N/A	N/A	N/A	N/A	N/A	7,7%
Partial response	6,9%	8,9%	5,9%	8,3%	6,4%	5,7%	7,1%
No one home /no answers	1,7%	4,7%	8,2%	9,3 %	14,0%	22,7%	8,5%
Moved outside Canada	4,7%	7,5%	5,1%	1,9%	2,5%	2,7%	4,2%
No interview – Mental/physical health problem	N/A	N/A	2,4%	4,9%	2,2%	1,6%	1,7%
Other non-responses	9,0%	21,0%	9,3%	8,3%	6,6%	4,3%	10,1%

9.2.5 Item Refusal and “Don’t Know” Rates

9.2.5.1 Refusal and “Don’t Know” Rates by Item

Items rates have been calculated from the number of refusal, “Don’t Know” and valid values for each variable, sub-module and module in the questionnaire. Derived variables and variables with few valid responses are not presented. Valid values exclude responses coded “not applicable” or “not stated”. Table 9.J shows refusal and “Don’t Know” rates for cycle 7 modules.

Table 9.J: Refusal and “Don’t Know” Rates by Module

Module	Refusal rate	“Don’t Know” rate
Overall	0.54%	1.05%
Admin	0.06%	0.07%
Household Record Variables (DHC)	0.09%	0.06%
General Health (GHC)	0.19%	0.25%
Sleep (SLC)	0.18%	0.86%
Height/Weight (HWC)	0.18%	0.17%
Nutrition (NU_, FV_, SK_, MK_)	0.39%	1.14%
Preventive Health (PHC)	0.49%	1.03%
Health Care Utilization (HCC)	0.22%	0.38%
Restriction of Activities (RAC)	0.36%	0.88%
Chronic Conditions (CCC)	0.04%	0.81%
Health Status (HSC)	0.19%	0.46%
Physical Activities (PAC)	1.49%	1.40%
UV Exposure (TUC)	1.26%	0.84%
Repetitive Strain (RPC)	0.19%	1.38%
Injuries (IJC)	0.15%	0.39%
Stress (STC, ST_)	0.25%	2.33%
Medication Use (DGC)	0.10%	0.55%
Smoking (SMC)	0.57%	1.38%
Alcohol (ALC)	0.29%	0.83%
Mental Health (MHC)	0.58%	0.74%
Social Support (SSC)	0.19%	0.72%
Socio-Demographic (SDC) - Language	0.77%	0.41%
Education (EDC)	0.44%	0.38%
Labour Force (LSC)	0.55%	0.73%

Module	Refusal rate	“Don’t Know” rate
Income (INC)	1.88%	3.10%
Food insecurity (FI_)	0.64%	0.51%

This table shows that refusal rates by module are very low and vary between 0.04% and 1.88%. The overall refusal rate is 0.54%. It tends to be the same variables or modules that have relatively high refusal rates in each cycle. Like the previous cycle, the income module has again the highest rate (1.88%). The physical activity module has the second highest refusal rate at 1.49%. The food insecurity modules, which was not included in the last three cycles, has a refusal rate of 0.64%, which is slightly higher than the overall rate. If we do more in depth analysis, although not shown in the table, rates of some labour force sub-modules are among the highest refusal rates in cycles 2 to 7, reaching 1.03% in cycle 7. In general, refusal rates by variable fluctuate between 0.0% and 2.11%, with the exception of some variables related to income and to physical activity that can reach refusal rates of 3.48% and 2.29% respectively.

Module “Don’t Know” rates are low and vary between 0.06% and 3.10%. The overall “Don’t Know” rate is 1.05%. It tends to be the same variables and modules that have relatively high “Don’t Know” rates in each cycle. For example, the income module has a “Don’t Know” rate of 3.10% and some sub-modules related to labour force and chronic condition have “Don’t Know” rates of 2.64% and 3.86%, which are among the highest rates in each cycle. The stress module has the second highest rate at 2.33% and three of four of its sub-modules have rates greater than 2.5%. In particular, the childhood and adult stressors, which was last included in cycle 4, has a 2.71% “Don’t Know” rate for cycle 7. Generally, the “Don’t Know” rates by variable fluctuate between 0.0% and 21.8% and income variables have the highest rates.

9.2.5.2 Refusal and “Don’t Know” Rates by Respondent

Refusal and “Don’t Know” rates were also calculated at the respondent level to determine the percentage of questions an individual refuses to answer or answers “Don’t Know”. All of these rates are for respondents who have a “completed” status for the cycle 7 interview. Table 9.K shows the proportion of respondents who did not refuse to answer any questions, who refused less than 1% and less than 3% of the questions asked. Also in this table are the rates for “Don’t know” answers to none of the questions, less than 0.5%, less than 1% and less than 5% of the questions asked.

Table 9.K: Refusal and “Don’t Know” Rates by Respondent

	Refusals to 0% of questions	Refusals to less than 1% of questions	Refusals to less than 3% of questions	“Don’t Know” to 0% of questions	“Don’t Know” to less than 0.5% of questions	“Don’t Know” to less than 1% of questions	“Don’t Know” to less than 5% of questions
Overall	90.2%	95.7%	96.3%	48.6%	70.2%	86.0%	95.8%
Males	90.2%	95.6%	96.1%	50.4%	70.9%	86.4%	95.7%
Females	90.2%	95.9%	96.5%	47.1%	69.6%	85.7%	95.9%
Under 12	100.0%	0.0%	0.0%	92.9%	92.9%	100.0%	0.0%
12-24	93.2%	96.9%	97.1%	43.3%	68.1%	84.1%	94.1%
25-44	92.9%	96.1%	96.3%	60.9%	79.8%	91.6%	96.6%
45-64	89.9%	96.1%	96.7%	51.2%	74.1%	89.6%	97.2%
65+	84.5%	93.8%	95.1%	29.9%	51.1%	73.3%	93.4%
Proxy	88.3%	94.9%	95.9%	39.8%	51.1%	69.7%	94.9%
Non-Proxy	90.3%	95.8%	96.4%	49.1%	71.2%	86.9%	95.8%

It can be seen that “Don’t Know” rates have more variation than the refusal rates when they are separated by sex, age group and interview type. For Cycle 7, 90.2% answered all of the questions and a little over 96% have refusal rates for less than 3% of the questions. As for “Don’t Know” rates, there are 48.6% of the respondents that answered all of the questions, in other words almost half of the respondents never answered “don’t know”. Around 70% of the respondents answered “Don’t Know” to less than 0.5% of the questions. When we look at 1% and 5% of the questions asked, the “don’t know” rates are 86% and 96% respectively. This shows that almost everyone who refuses or responds “Don’t Know” does so for only a few questions.

10. Guidelines for Tabulation, Analysis and Release

This section of the documentation outlines the guidelines that should be followed by users to tabulate, analyze, release or otherwise publish any data derived from the NPHS data. With the aid of these guidelines, users 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.

10.1 Rounding Guidelines

In order that dissemination of estimates derived from NPHS data corresponds to estimates produced by Statistics Canada, Users should use the following guidelines regarding the rounding of such estimates. Unrounded estimates imply greater precision than actually exists.

- a) Estimates in the main body of a statistical table should 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 should 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 should be computed from unrounded components (i.e., numerators and/or denominators) and then, they are to be rounded to one decimal using normal rounding. In normal rounding to a single digit, 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 increased by 1.
- d) Sums and differences of aggregates (or ratios) should 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, it is suggested to users to note the reason for such differences in the publication or release document(s).

10.2 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 statistical tables, users must apply the proper

sampling weight. If proper weights are not used, the estimates derived from the various subsets of respondents cannot be considered representative of the 1994/1995 target population, and will not correspond to those produced by Statistics Canada.

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

10.2.1 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 computed.

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 quit smoking between cycles 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:

At the present, do/does ... smoke cigarettes daily, occasionally or not at all? (SMCB_2)

- 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 increase in the number of cigarettes smoked per day by daily smokers who had an increase in consumption between two cycles.

Example of Quantitative Question:

How many cigarettes do/does you/he/she smoke each day now? (SMCB_4)

|_|_| Number of cigarettes

10.2.2 Tabulation of Categorical Estimates

Estimates of the number of people with a certain characteristic can be obtained 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) by summing the final weights of records having the characteristic of interest for the numerator (\hat{X}),
- b) by summing the final weights of records having the characteristic of interest for the denominator (\hat{Y}), then
- c) by dividing the numerator estimate by the denominator estimate.

10.2.3 Tabulation of Quantitative Estimates

Estimates of sums or averages for quantitative variables can be obtained using the following three steps (only step a) is necessary to obtain the estimate of a sum):

- 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 the estimate of the average number of cigarettes smoked each day by individuals who smoke daily, first compute the numerator (\hat{X}) by summing the product between the value of variable **SMCB_4** and the final weight. Next, sum this value over those records with a value of "daily" to the variable **SMCB_2**. The denominator (\hat{Y}) is obtained by summing the final weight of those records with a value of "daily" to the variable **SMCB_2**. Divide (\hat{X}) by (\hat{Y}) to obtain the average number of cigarettes smoked each day by daily smokers.

10.3 Guidelines for Statistical Analysis

The NPHS is based upon a complex sampling 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 final weight divided by the average of the final weights for the sampled units (people) contributing to the estimate in question.

CV tables were produced in the past for the cross-sectional data. CV tables were not created for the longitudinal files, as a very large number of possible variable combinations for analysis exist. To correctly estimate the variance, NPHS recommends the use of the bootstrap method. With the bootstrap method, the complexity of the weighting and the survey design are incorporated into the calculation of the variance. A SAS bootstrap variance program, along with accompanying documentation and examples of how to use it, has been created to facilitate the calculation of the variance using the bootstrap method. The program also calculates the accompanying coefficient of variation. A similar version of the program is also available in SPSS. It is important for users to learn how to use it as the program will generate exact estimates of individual variances to assess the quality of tabulated estimates and is highly recommended over the use of the scaled weights approach. Some statistical packages such as STATA have the ability to read in the stratum and cluster information to use in variance estimation, which improves the quality of the estimate but does not take into account the different adjustments applied to the weights.

10.4 Release Guidelines

Before releasing or publishing any total or proportion estimates from the master files, users must first determine the number of sampled respondents having the characteristic of interest (for example, the number of respondents who smoke when interested in the proportion of smokers for a given population). If this number is less than 10, the weighted estimate should not be released regardless of the value of the coefficient of variation for this estimate. This is due to the fact that the possibility of obtaining an artificially low variance is greater with a sample size less than 10. For weighted estimates based on sample sizes of 10 or more, users should determine the coefficient of variation of the estimate and follow the guidelines described in Table 10.A.

Table 10.A: Sampling Variability Guideline

Type of Estimate	C.V. (in %)	Guidelines
Acceptable	0.0 - 16.5	Estimates can be considered for general unrestricted release. Requires no special notation.
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 E (or in some other similar fashion).
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 F (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>

11. Using the Longitudinal Master File

11.1 Use of Longitudinal Weights

The Cycle 7 master file contains all 17,276 panel members (including occasional or permanent non-respondents). Four subsets of respondents (Section 7.7) to which correspond a distinct set of weights (Section 8.1) have been created. Flags were created to identify records that are part of a particular subset (Table 8.A). Records that are not part of a particular subset have a flag equal to zero and the weight variable set to blank for that particular subset. **To create the subset of interest, select those records that have the appropriate flag variable equal to one.**

Weight WT64LS is called the “square weight” and applies to the 17,276 members that make up the original longitudinal sample. All non-response should be taken into account for any calculation.

Weight WT6BLF is called the “Longitudinal Full” weight and applies to the 10,992 records that are included in the “Full” subset of respondents.

Weight WT6BSLS is called the “Longitudinal Square Share” weight and applies to the 16,007 records that are included in the “Square Share” subset of respondents.

Weight WT6BSLF is called “Longitudinal Full Share” weight and applies to the 10,668 respondents that are included in the “Full Share” subset of respondents.

11.2 Ensuring the Reliability of Estimates with the Use of Bootstrap Weights

Bootstrap weights are necessary for variance estimation. Information on the bootstrap method for variance estimation can be found in Section 9.1.1. Each subset of respondents has a set of bootstrap weights associated with it. Four different sets of bootstrap weights were created for the Cycle 7 data: the square, the square share, the full and the full share. For more information on these subsets, see Section 7.7. Table 11.A presents the subset of respondents with their corresponding bootstrap file name.

Table 11.A: Subsets of Respondents and Corresponding Bootstrap Weights files

Subset of respondents	Name of the Bootstrap Weights file	Number of Respondents
Longitudinal Square	B5long	17,276
Longitudinal Full	B5lngf	10,992
Longitudinal Square Share	B5long_share	16,007
Longitudinal Full Share	B5lngf_share	10,668

Due to the complex sample design, users should use the supplied Bootvar program for variance calculation. The standard variance output from other statistical packages such as SAS and SPSS may grossly underestimate the variance of an estimate for this survey. **It is**

the responsibility of the user to ensure the quality/reliability of the estimates that they are producing by following the guidelines laid out in Chapter 10 and correctly calculating the variance for all estimates. Failure to do so could lead to some misinterpretation of results and jeopardize the quality of the research work.

Some statistical software are capable of including the stratum and cluster information as input when performing analytical processing, which does provide a variance estimate much closer to the true variance estimate, but these packages fail to account for the various weighting adjustments, which in some cases can impact the variance estimates considerably.

11.3 Variable Naming Convention

NPHS has adopted a variable naming convention that 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 (1994/1995, 1996/1997, 1998/1999, 2000/2001, 2002/2003, 2004/2005, and 2006/2007) 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 1994/1995, 1996/1997, 1998/1999, 2000/2001, 2002/2003, 2004/2005, and 2006/2007, and the variable names should only differ in the position that identifies the particular survey occasion in which they were collected. This convention is followed throughout the longitudinal survey, and is adopted by all NPHS surveys: the household component, the health institutions component, and previously the North component and supplements.

11.3.1 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 name / Questionnaire section name
Position 3:	Survey type
Position 4:	Cycle
Position 5:	Variable type
Positions 6-8:	Variable number / name from questionnaire

For example: the variables DHC4_AGE, DHC6_AGE, DHC8_AGE, DHC0_AGE, DHC2_AGE, DHCA_AGE and DHCB_AGE:

DH: in the Demographic and Household content section of the questionnaire;

C: questions which are Core content on the household survey;

4/6/8/0/2/A/B: Cycle 1 (1994/1995) variable / Cycle 2 (1996/1997) variable / Cycle 3 (1998/1999) variable / Cycle 4 (2000/2001) variable / Cycle 5 (2002/2003) variable / Cycle 6 (2004/2005) variable / Cycle 7 (2006/2007) variable;

_: Cycle specific focus content;

AGE: the variable name.

11.3.2 Positions 1-2: Variable Name / Questionnaire Section Name

AD	Alcohol dependence
AL	Alcohol Consumption
AM	Administration (of the survey)
AP	Attitude toward parents
BF	Breast-feeding
BP	Blood pressure
CC	Chronic conditions
CO	Coping – Stress - Alberta buy-in, Cycles 1 and 2
CE	Contact/Exit
DG	Drug/medication use
DH	Demographics and household variables
DV	Dental visits
ED	Education
ES	Emergency services
EX	Eye examination
FH	Personal and family medical history
FI	Food insecurity - HRDC buy-in, Cycle 3
FS	Flu shots
FV	Fruit and vegetable consumption
GE	Geographic identifiers (methodology)
GH	General health
HC	Health care utilisation
HI	Health information
HH	Household
HN	Health number
HS	Health status
HV	HIV
HW	Height, weight and body image
IJ	Injuries
IN	Income
IS	Insurance
LF/LS	Labour force
MH	Mental health
MK	Milk consumption

NU Nutrition
Positions 1-2: Variable Name / Questionnaire Section Name (continue)

PA	Physical activities
PC	Physical check-up
PH	Preventive health
PR	Province
PY	Psychological resources (self-esteem, mastery, sense of coherence)
RA	Restriction of activity
RH/MB	Residential history
RP	Repetitive strain
RS	Road safety
RT	Rationality - Manitoba buy-in, Cycle 1
SC	Self-care
SD	Socio-demographic characteristics
SH	Sexual health
SK	Soft drink consumption
SL	Sleep
SM	Smoking
SP	Sample identifiers (methodology)
SS	Social support
ST	Stress
SV	Health care utilisation
TA	Tobacco alternatives (Health promotion 1998)
TU	Tanning -UV exposure
TW	Two-week disability
VS	Violence and personal safety
WF	Subset flags
WT	Sample weights (methodology)

A few important identifying variables do not follow the naming convention: e.g. REALUKEY, PERSONID, CYCLE, SUBCYCLE, DESIGPRV, STRATUM, and REPLICAT.

There are also some variables that are considered “constant”. Table 11.B presents the variables that appear only once of the data file. The name of these variables does not follow the naming convention.

Table 11.B: “Constant” Longitudinal Variables

Longitudinal Variable Name	Concept
AOI	Age at Time of Immigration
COB	Country of Birth
COBC	Code for Country of Birth
COBGC	Code for Country of Birth (7 groups) - Grouped
COD10	Cause of Death Code (ICD-10)
DESIGPRV	Province of Residence in 1994
DOB	Date of Birth

Longitudinal Variable Name	Concept
DOD	Day of Death
HWB	Birth Weight
HWBG1	Birth Weight - Grouped
IMM	Immigration Status
MOB	Month of Birth
MOD	Month of Death
REPLICATE	Replicate
SEX	Sex
STRATUM	Stratum
YOB	Year of Birth
YOD	Year of Death
YOI	Year of Immigration to Canada

11.3.3 Position 3: Survey Type

A	Asthma supplement
B	Province-specific buy-in content – children’s questions
C	Core questions repeated in each cycle
F	Food Insecurity supplement
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 focus questions, not repeated in every cycle (e.g., food choice in Cycle 3, 5 and 7)
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) (for example, weights, agreement to share, date of interview variables, etc.)

11.3.4 Position 4: Cycle (years)

4	Cycle 1 (1994/1995)
6	Cycle 2 (1996/1997)
8	Cycle 3 (1998/1999)
0	Cycle 4 (2000/2001)
2	Cycle 5 (2002/2003)
A	Cycle 6 (2004/2005)
B	Cycle 7 (2006/2007)
C	Cycle 8 (2008/2009)
D	Cycle 9 (2010/2011)
E	Cycle 10 (2012/2013)

11.3.5 Position 5: Variable Type

Position 5	Variable type	Description
–	Collected variable	A variable that appeared directly on the questionnaire
C	Coded variable	A variable coded from one or more collected variables (e.g., North American Industry Classification System (NAICS))
D	Derived variable	A variable calculated from one or more collected or coded variables, usually calculated during head office processing (e.g., Comprehensive Health Status Measurement System (CHSMS-HUI3))
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). It can also denote that a long answer was collected (e.g., restriction of activity 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

11.3.6 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 DHC4_AGE in Cycle 1 (1994/1995), DHC6_AGE in Cycle 2 (1996/1997), etc. 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 STC2_C1; the first question on work stress was named STC2_W1. 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. Access to NPHS Data

12.1 Research Data Centres

Confidentiality concerns preclude general dissemination of longitudinal NPHS data in public use microdata file (PUMF) format. However, access to all the longitudinal master microdata files including the data for cycles 1 to 7 is available through Statistics Canada's Research Data Centres (RDCs) program. The RDCs program is part of an initiative by Statistics Canada, the Social Sciences and Humanities Research Council (SSHRC) and university consortia to help strengthen Canada's social research capacity and to support the policy research community.

RDCs provide researchers with access, in a secure university setting, to microdata from population and household surveys. The centres are staffed by Statistics Canada employees. They are operated under the provisions of the *Statistics Act* in accordance with all the confidentiality rules and are accessible only to researchers with approved projects who have been sworn in as "deemed employees". RDCs are located throughout the country, so researchers do not need to travel to Ottawa to access Statistics Canada microdata. More information is available at the Research Data Centre Program web site: <http://www.statcan.ca/english/rdc/index.htm>.

12.2 Remote Access

A second option, if the RDCs are not accessible for the researcher, is Health Statistics Division's Remote Access service. This service provides researchers with a means to develop and test their own computer programs using synthetic files that mimic the actual master files. Researchers then submit their programs to a dedicated e-mail address. The programs are run against the master microdata files on an internal secure server, outputs are vetted for confidentiality, and sent back to the researcher by return e-mail. For more information on this service, please contact the Data Access and Information Services team at nphs-ensp@statcan.ca.

12.3 Data Liberation Initiative

PUMFs are available for each of the first three cycles of the NPHS, providing widespread access to the cross-sectional components of the survey. The NPHS PUMFs can be accessed through the Data Liberation Initiative (DLI) at participating Canadian universities and colleges. For more information, please consult the following Statistics Canada's link: <http://www.statcan.ca/english/Dli/dli.htm>. Cycles 1, 2 and 3 NPHS PUMFs can also be purchased. To this end, please contact the Data Access and Information Services team at hd-ds@statcan.ca or one of Statistics Canada's Regional Offices.

12.4 Analytical Reports and Tabulations

With the release of NPHS Cycle 5 data, results from the survey were presented in a free Internet Publication entitled "*Healthy Today, Healthy Tomorrow? Findings from the National Population Health Survey*". The publication (catalogue 82-618M) is centered on a series of articles addressing important health issues using NPHS longitudinal data. To consult this publication, use the following link:

<http://www.statcan.ca/bsolc/english/bsolc?catno=82-618-M&CHROPG=1>

Longitudinal Cansim tables are also available free of charge on the Statistics Canada Internet site. They present changes, from one NPHS cycle to another one, in smoking, self-rated health, body mass index and physical activity. One can access Cansim tables by clicking the above publication link, and then chose the Data tables options on the left side of the main page.

Research articles based on the NPHS often appear in *Health Reports*, a quarterly journal produced by Health Statistics Division. This product is available as a standard printed publication (catalogue no. 82-003-XPE) or in electronic format (catalogue no. 82-003-XIE) on the Statistics Canada Internet site as. To obtain more information, consult the following links:

<http://www.statcan.ca/bsolc/english/bsolc?catno=82-003-X&CHROPG=1>

<http://www.statcan.ca/bsolc/english/bsolc?catno=82-003-S&CHROPG=1>

Custom tabulations based on the NPHS cross-sectional data are also available on a cost recovery basis. No custom tabulation is produced with the NPHS longitudinal data. For more information or estimates on costs and feasibility, contact the Health Statistics Division's Data Access and Information Services at hd-ds@statcan.ca.

13. An Analytical Technique for Longitudinal Survey Data

Longitudinal surveys like the NPHS have the advantage of watching their analytical potential grow over time. These longitudinal surveys allow health analysts to study events that affect health in the life of an individual, to study their effects and causes as well as to produce incidence rates. The analysis of such data equally highlights the complexity of the network of relations that exist between the health of an individual and the intensity of exposures to different risk factors. This advantage permits research of more complex research questions about health. On the other hand, certain issues reduce this analytical potential and may even pose a long term risk to longitudinal surveys.

The analysis of data from longitudinal surveys faces different issues than those from cross-sectional surveys. On the one hand, studying more complex research questions leads data users to use more complex analytical techniques, often those that are less well-known and less well-documented. On the other hand, attrition of data over time (see section 9.2.4) is one of the most important challenges in the analysis of longitudinal data. In addition to the risk of introducing a bias in the estimates, a decrease in sample size due to attrition can also prove to be problematic not only for the variety of statistical analyses but also for their quality.

In addition, the most common and simplest analytical approach is to use the longitudinal full subset (see sections 7.6 and 7.7). This subset, however, is the most sensitive to attrition, which reduce the number of individuals in every cycle for this subset. After seven cycles, the cumulative attrition rate for the NPHS Full subset has reached 36.4% In other words, 63.6% of the original panel of respondents has a complete response (see section 9.2.4) to the seven cycles of the NPHS.

The purpose of this chapter is to give users a practical analytical tool that will enable them to reduce the effects of attrition on their analyses and thereby extend the analytical potential of the data. Consequently, mitigating the effects of attrition will increase the accuracy of the estimates and the relevance of the survey

13.1 Cycle Twinning Approach

The cycle twinning approach (also known as “Pooling of Repeated Observations”) is an analytical approach that allows the effects of erosion to be reduced. It consists of using a subset of cycles from a respondent as the unit of analysis rather than considering all cycles from a respondent as the unit of analysis. This approach is a particular case in a type of statistical models called “marginal models”²³.

Although the fictitious example and the application described in this section are limited to two consecutive cycles completed by the same respondent as the unit of analysis, this approach can be adapted to many different scenarios. For example, the length of exposure to different risk factors could be shorter or longer, the number of cycles (or repeated measures) could be 1 or more than 2, the dependent variables could be continuous or discrete with temporary or permanent conditions such as death or certain chronic conditions..

²³ Fitzmaurice, G.M., Laird, N.M., and Ware, J.H. (2004). *Applied Longitudinal Analysis*. New York: Wiley

Fictitious example

Thus after seven cycles of collection, when the Cycle twinning approach is used, a respondent can supply up to six units of analysis if he or she belongs to the Full subset : Cycles 1-2, Cycles 2-3, Cycles 3-4, Cycles 4-5, Cycles 5-6 and Cycle 6-7(see respondent A in Table 13.A). If a respondent did not respond to cycle 3, 4 and 7 for example, but responded to all other cycles, he or she can produce two units of analysis: Cycles 1-2 and Cycles 5-6 (see respondent B in Table 13.A). The latter respondent is not part of the Cycle 7 Full subset given his non-response to cycles 3, 4 and 7 and he or she would be automatically excluded from any analyses using this subset. On the other hand, by using the cycle twinning approach, two analytical units can come from this respondent. In order to supply at least one analytical unit, all that is required is that a respondent answers two consecutive cycles. According to Table 13.A, respondents A, B, C and D will provide 6, 2, 4 and 1 analytical units(s) respectively.

Table 13.A: Example of Profiles of Response from Fictitious NPHS Respondents

Respondent	Cycles with a complete response						
	Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5	Cycle 6	Cycle 7
A	X	X	X	X	X	X	X
B	X	X			X	X	
C	X	X	X	X	X		X
D	X		X	X			X

According to this fictitious table, the Cycle 7 Full subset would contain only one of the four respondents since the analytical unit using the Full subset is respondents who answer all cycles. Thus, for this example, the attrition rate after seven cycles would be 75%.

When the cycle twinning approach is used, the analytical unit becomes a record that contains two consecutive cycles from one respondent. Thus, according to the example in Table 13.A, the potential number of analytical units is 24 (6 possible combinations of two consecutive cycles multiplied by four respondents). As calculated above, the total number of analytical units is 12 (6+2+4+1), in other words 50% of the potential maximum number of analytical units. Consequently, the attrition rate decreases to 50% compared to 75% for the Full subset.

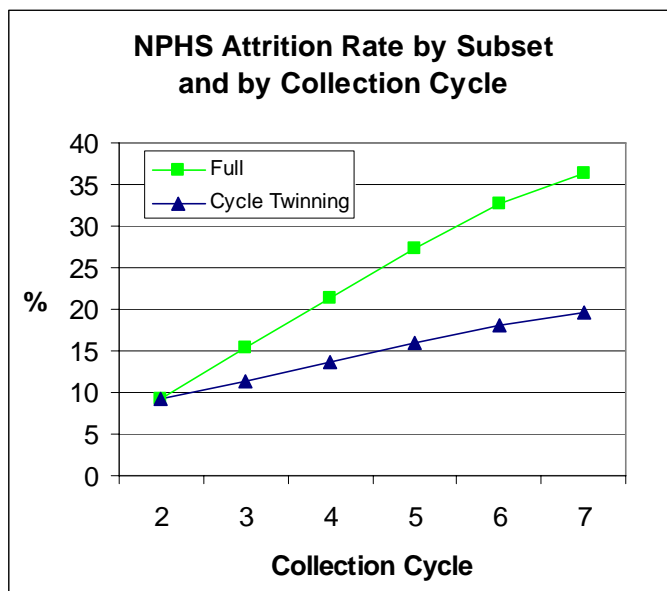
By applying this approach to the NPHS, the maximum potential sample after seven cycles of collection reaches 103,656 (6 x 17,276) analytical units (the NPHS contains 17,276 respondents multiplied by 6 analytical units possible after seven cycles). After seven cycles of collection, attrition reduces the number of potential analytical units to 83,321, in other words, to 80% of the original sample resulting in an attrition rate of 20%.

In comparison, the attrition rate for the Full subset is 36% after seven cycles. By taking the analytical units, that result from using the cycle twinning approach, the attrition rate diminishes by approximately 16 percentage points in comparison to that observed with the Full subset (see Table 13.B and Figure 13.A).

Table 13.B: Attrition Rates

Cycle	Subset	
	Full	Cycle twinning approach
2	9.3%	9.3%
3	15.3%	11.4%
4	21.3%	13.7%
5	27.3%	15.9%
6	32.7%	18.0%
7	36.4%	19.6%

Figure 13.A: NPHS Attrition Rates Taken from Table 13.B by Subset and by Collection Cycle



13.2 Creation of the Modified Subset

Before being able to analyze the data using the cycle twinning approach, it is necessary to modify the structure of the NPHS Square subset. The NPHS Square subset contains 17,276 records or analytical units where each record represents information from all cycles for a single respondent. The transformation of this subset consists of modifying the analytical unit so that each unit consists of two consecutive cycles from a single respondent. For Cycle 7, the number of records for this modified subset would be 83,321 (see section 13.1).

Using the example from Table 13.A, the record from respondent A is re-written to 6 analytical units, respondents B, C and D to 2, 4 and 1 analytical unit(s) respectively (see

Table 13.C). Here are the steps to follow for modification of the NPHS Square subset and for conducting analysis using the cycle twinning approach.

Steps to follow for the creation of the NPHS modified subset:

- Step 1: Calculate the possible number of twinings for two consecutive cycles for each of the 17,276 respondents.
- Step 2: Clone each NPHS respondent the same number of times as the number of twinings calculated in step 1.
- Step 3: Identify the longitudinal variables of interest for the twinned cycles.
- Step 4: Identify the sample weights from cycle 1: the longitudinal square weights (WT64LS), as well as the corresponding bootstrap weights (B5long) (see section 11.1 and 11.2).
- Step 5: Keep only the information resulting from steps 3, 4 and 5. This step will increase the efficiency of analytical programs.

Table 13.C: Example of the Modified Subset Structure for Twinning

Respondent	Twinned Cycles	Cycle 1 weight	Bootstrap weights	Variables of interest
A	1-2	W_A	BS1 _A à BS500 _A	
A	2-3	W_A	BS1 _A à BS500 _A	
A	3-4	W_A	BS1 _A à BS500 _A	
A	4-5	W_A	BS1 _A à BS500 _A	
A	5-6	W_A	BS1 _A à BS500 _A	
A	6-7	W_A	BS1 _A à BS500 _A	
B	1-2	W_B	BS1 _B à BS500 _B	
B	5-6	W_B	BS1 _B à BS500 _B	
C	1-2	W_C	BS1 _C à BS500 _C	
C	2-3	W_C	BS1 _C à BS500 _C	
C	3-4	W_C	BS1 _C à BS500 _C	
C	4-5	W_C	BS1 _C à BS500 _C	
D	3-4	W_D	BS1 _D à BS500 _D	

13.3 Methodological Aspects of the Cycle Twinning Approach

As mentioned in section 10.3, the sampling design and the sampling probabilities have an impact on the estimation method and variances calculations. The cycle twinning approach brings with it an additional layer of complexity: analytical units from the same respondent are highly correlated to each other.

With the help of “marginal models” (Fitzmaurice et al. (2004)), it is enough to consider the analytical unit resulting from the cycle twinning approach as an additional stage in the

NPHS sampling design in order to account for the dependency between the analytical units in the variance calculation. The correlation between the analytical units from the same respondent is therefore included in the sample design effect. Creating the modified subset as described above and using the bootstrap method for variance calculation will ascertain the analyst to compute the variance correctly.

The choice of the variance calculation method is even more important given the high correlation between analytical units from the same respondent. The use of one of the methods described in section 10.3, like the bootstrap method (SAS and SPSS) or other software such as STATA or SUDAAN is highly recommended.

13.4 An Example of How to Use the Cycle Twinning Approach: Quitting Smoking

Quitting smoking is one of the most important steps that smokers can take to improve their health. Understanding the factors that are associated with smoking cessation is important for public health programs aimed at reducing the smoking rate. Shields (2005)²⁴ used the NPHS to analyze factors associated with smoking cessation among people aged 18 or older who were daily smokers.

For this analysis, cycles 1 to 5 of the NPHS were used. The NPHS follows the same sample of individuals over time, interviewing them at two-year intervals. An analysis file was assembled by examining successive pairs of NPHS cycles. The target population of interest was all respondents aged 18 or older who reported that they were daily smokers. Daily smokers were considered to be quitters if, in the following cycle, they reported not smoking at all in the last two years.

The analysis file was created by identifying daily smokers aged 18 or older at four baseline cycles (1, 2, 3, and 4) and determining if they had quit at the follow-up cycle two years later (i.e., cycles 2, 3, 4 and 5, respectively). For a record to be written to the analyses file, two conditions had to be satisfied:

- The individual was a daily smoker aged 18 or older at the baseline cycle²⁵.
- Smoking status was known at the follow-up cycle two years later (i.e., it was known if the individual still smoked or had quit),

The analysis file used contains many individuals who contributed more than one record, because they reported that they were daily smokers in more than one survey cycle.

To illustrate, an artificial example is shown below. Suppose that only two risk factors are studied to model the probability of quitting smoking: the number of cigarettes smoked per day (x) and the age of smoking initiation (z). Let y denote a binary variable that takes on the value “1” if a daily smoker had stopped smoking in the subsequent cycle two years later; otherwise y has the value “0”. Information about x , y , and z is collected at each cycle.

Let p denote the probability that y takes the value “1”; i.e., the individual has stopped smoking. Suppose that the model of interest is the logistic model:

²⁴ Shields, M. (2005). The journey to quitting smoking. Health Reports, Vol. 16, No. 3. Statistics Canada, Catalogue 82-003.

²⁵ In a pair of twinned-cycles the baseline cycles refers to the first cycle of the pair. For example, for twinned cycles 3 and 4 the baseline cycle is Cycle 3.

$$\log\left(\frac{p}{1-p}\right) = \alpha_0 + \alpha_1 x + \alpha_2 z,$$

and that the objective of the study is to estimate α_0 , α_1 and α_2 and standard errors of their estimates. Suppose that only the first four cycles of the NPHS are being used for this analysis, and there are only four individuals on the NPHS file (all men) resulting in the following analysis file:

Cycle	Person ID	x	Z	Y ²⁶
1	1	x_{11}	z_{11}	0
1	2	x_{21}	z_{21}	1
1	3	x_{31}	z_{31}	1
1	4	x_{41}	z_{41}	0
2	1	x_{12}	z_{12}	0
2	4	x_{42}	z_{42}	0
3	1	x_{13}	z_{13}	0
3	3	x_{33}	z_{33}	1
3	4	x_{43}	z_{43}	1

- Person 1, who was a daily smoker in all 4 cycles, contributes 3 records to the analysis file since he was 18 or older and a daily smoker for all three baseline cycles (1, 2 and 3).
- Person 2 contributes only one record. He was a daily smoker at cycle 1, had quit at cycle 2, and did not respond to the survey at cycles 3 and 4.
- Person 3 contributes two records. He was a daily smoker in cycle 1, quit smoking in cycle 2, took up smoking again in cycle 3, and quit by cycle 4.
- Person 4 contributes three records. He was a daily smoker in the first 3 cycles and then quit in cycle 4.

This pattern of responses can be summarized as follows:

Person ID	Cycle 1	Cycle 2	Cycle 3
1	X	X	X
2	X		
3	X		X
4	X	X	X

²⁶ Please note this variable is derived from the subsequent cycle two years later.

Suppose that the analysis is conducted using a software package for design-based logistic regression that requires an input data file that includes both the survey weight and the corresponding bootstrap weights. The survey and bootstrap weights from cycle 1 are used. The analysis file must be merged with the bootstrap weight for cycle 1, and should have the following structure:

Record	Person	Wt	X	Z	Y	bs1	bs2	---
1	1	w_1	x_{11}	z_{11}	0	$bs1_1$	$bs2_1$	---
2	1	w_1	x_{12}	z_{12}	0	$bs1_1$	$bs2_1$	---
3	1	w_1	x_{13}	z_{13}	0	$bs1_1$	$bs2_1$	---
4	2	w_2	x_{21}	z_{21}	1	$bs1_2$	$bs2_2$	---
5	3	w_3	x_{31}	z_{31}	1	$bs1_3$	$bs2_3$	---
6	3	w_3	x_{33}	z_{33}	1	$bs1_3$	$bs2_3$	---
7	4	w_4	x_{41}	z_{41}	0	$bs1_4$	$bs2_4$	---
8	4	w_4	x_{42}	z_{42}	0	$bs1_4$	$bs2_4$	---
9	4	w_4	x_{43}	z_{43}	1	$bs1_4$	$bs2_4$	---

Where w_i is the survey weight variable for person i and bsj_i is the j^{th} corresponding bootstrap weight for the same person. Note that the weights and bootstrap weights are the same for multiple observations from the same person.

This approach uses the same design information (weights and bootstrap weights) and the same software that would be used if the model of interest were being fit to only one observation per individual. The approach simply requires that an individual's weight and bootstrap weights be assigned to multiple observations for that individual.

The question about which weight and bootstrap weights variables are the most appropriate can be debated, but the weight variables from the first cycle of the survey are often a good choice. This allows the analyst to maximize the sample size by making use of partial information when an individual has not responded to all NPHS cycles.

Although additional complexities are associated with partial response within individuals, these are not discussed in this document.

Appendix A: NPHS Household Component, Changes to the Questionnaire for Cycle 7 (2006/2007)

1. Introduction

This appendix describes the changes between the Cycle 6 (2004/2005) and the Cycle 7 (2006/2007) questionnaires. Some questions from Cycle 6 were removed in Cycle 7 and other questions, not in Cycle 6, were added in Cycle 7. The changes between the two cycles are described in detail in point 3. Globally, the main additions to Cycle 7 are in the Nutrition module (12 focus questions on food choice were brought back from Cycle 5), in the Stress module (7 focus questions on childhood and adult stressors (trauma) were brought back from Cycle 4) and the 3 questions from the Food insecurity module were also brought back from Cycle 3. Mainly, the questions removed at Cycle 7 are 16 focus questions about coping from the Stress module.

2. Changes to Questionnaire Structure

Apart from the questions added and removed at Cycle 7, the order of the questionnaire remained the same as in Cycle 6.

3. Changes to Core Content

In the following description, external question names from the questionnaires are used. Some internal question names may have been changed to ensure consistency throughout the questionnaire due to the deletion or addition of questions. The variable names for the master, share, and public files are created using the variable naming convention.

Sections without modifications

- Household Record Variables
- General Health (GH)
- Sleep (SL)
- Height and Weight (HW)
 - Body Image
- Nutrition (NU)
 - Supplement use
 - Fruit and vegetable consumption
 - Soft drink consumption
 - Milk consumption
- Preventive Health (PH)
- Health Care Utilization (HC)
 - Home Care
- Restrictions of Activities (RA)
- Chronic conditions (CC)
 - Food or Digestive Allergies
 - Other Allergies
 - Asthma
 - Fibromyalgia
 - Arthritis or Rheumatism excluding Fibromyalgia

- Back Problems
- High Blood Pressure
- Migraine Headaches
- Chronic Bronchitis or Emphysema
- Diabetes
- Epilepsy
- Heart Disease
- Cancer
- Intestinal or Stomach Ulcers
- Effects of a stroke
- Urinary Incontinence
- Bowel Disorder
- Alzheimer’s Disease or other Dementia
- Cataracts
- Glaucoma
- Thyroid Condition
- Other Long-Term Condition
- Health Status (HS)
 - Vision
 - Hearing
 - Speech
 - Getting Around
 - Hands and Fingers
 - Feelings
 - Memory
 - Thinking
 - Pain and Discomfort
- Physical Activities (PA)
- UV Exposure (UV) (TU)
- Repetitive Strain (RP)
- Injuries (IJ)
- Stress (ST)
 - Ongoing Problems
 - Work Stress
 - Mastery
- Medication Use (DG)
- Smoking (SM)
- Alcohol (AL)
- Mental Health (MH)
- Social Support (SS)
- Language (SD)
- Education (ED)
- Labour Force (LF)
 - Job Attachment
 - Job Search – Last 4 Weeks
 - Past Job Attachment
 - Job Description
 - Absence/Hours

- Other Job
- Weeks Worked
- Looking for Work
- Income (IN)
- Provincial Health Number and Administration (AM)
 - Provincial Health Number
 - Administration

Sections with modifications

Nutrition (NU)

- Food choice
 - 12 Focus questions (NU_B_1A, NU_B_1C to NU_B_1E, NU_B_2A to NU_B_2C, NU_B_3A to NU_B_3D and NU_2_3G) from this sub-section of Cycle 5 were added in Cycle 7.

Stress (ST)

- Childhood and adult stressors:
 - 7 Focus questions (ST_B_T1 to ST_B_T7). Those questions were brought back from Cycle 4 to catch respondents who turned 18 years old since Cycle 4.
 - Coping (CO):16 Focus questions (CO_A_1 to CO_A_16) from Cycle 6 were dropped.

Food insecurity (FI)

- 3 Focus questions (FI_B_1 to FI_B_3) were brought back from Cycle 3 to be able to do comparisons between cycles.

4. Buy-in Content

No buy-in content in Cycle 7.

Appendix B: NPHS Household Component, Examples of Variables from Previous Interviews Used as Additional Information in Cycle 7 (2006/2007)

Blood Pressure; Mammography; Pap Smear Test	In Cycle 1 and Cycle 2 the respondent was asked whether he/she ever had his or her blood pressure taken (or ever had a mammography etc.). In Cycle 3 the questions were repeated; however, the respondent was probed when said that he or she has not had the test done and in the previous cycle reported the contrary. In Cycles 4, 5, 6 and 7, if the respondent had reported that he or she had had the test performed in a previous interview, only the question on the last time it was done was asked.
Restriction of Activities	Information on whether or not the respondent had a disability in a previous interview was used. If the status changed, an explanation of that change was probed.
Chronic Conditions	For each respondent, response to selected chronic conditions (asthma, fibromyalgia, arthritis, high blood pressure, migraine headaches, diabetes, epilepsy, stomach or intestinal ulcers and the effects of a stroke) in a previous interview were used to help explain change. If it was a newly acquired condition, the date of onset for the condition was captured.
Smoking	If a daily smoker had reported the age at which he/she started smoking daily during last interview, that response was used in Cycle 7. For the occasional smoker or non-smoker in Cycle 7 who had reported smoking daily (or having ever smoked daily) during last interview, a flag about daily smoking was re-input. If smoking status changed, an explanation of that change was probed.
Socio-demographic Characteristics	For all respondents, a flag indicating that country of birth had been collected was used. 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 of the main job from the previous interview were fed back. If the respondent indicated that they worked in the previous year, they were asked to confirm the employer name. When there was a change, the information was collected
Health Number	A flag indicated whether the health number that was collected in an earlier interview was valid. If the respondent's health number had not changed since last cycle or was invalid then the health number was asked again.