
The Adult Literacy and Life Skills Survey, 2003 and 2008

Public Use Microdata File

User's Manual

Table of Contents

1.0	Introduction.....	1
2.0	Survey Overview.....	2
3.0	Survey Objectives.....	3
4.0	Concepts and Definitions	4
4.1	Defining and Measuring Proficiency: Overview.....	4
4.2	Understand what was measured in the International Literacy and Skills Survey.....	5
4.2.1	Introduction	5
4.2.2	Scaling the literacy, numeracy and problem solving tasks in ALL	6
4.3	Measuring prose and document literacy in ALL	7
4.3.1	Identifying task characteristics.....	8
4.3.2	Type of match	10
4.3.3	Type of information requested	10
4.3.4	Plausibility of distractors	11
4.3.5	Characterizing prose literacy tasks.....	11
4.3.6	Characterizing document literacy tasks.....	15
4.4	Measuring numeracy in ALL	19
4.4.1	Identifying task characteristics.....	21
4.4.2	Everyday life	21
4.4.3	Work-related	21
4.4.4	Societal or community	21
4.4.5	Further learning	22
4.4.6	Quantity and number	23
4.4.7	Dimension and shape	23
4.4.8	Pattern, functions and relationships	24
4.4.9	Data and chance.....	24
4.4.10	Change	24
4.4.11	Characterizing numeracy tasks	25
4.5	Measuring problem solving in ALL	29
4.5.1	Identifying task characteristics.....	31
4.5.2	Characterizing problem solving tasks.....	32
4.6	Conclusion	36
4.6.1	Some analytical considerations	37
5.0	Survey Methodology	43
5.1	Assessment design.....	43
5.2	Target population and sample frame.....	44
5.3	Sample design	46
5.4	Sample size	50
5.5	Data collection	50
5.6	Scoring of tasks	54
5.7	Survey response and weighting	56
6.0	Survey Procedures and Data Processing	59
6.1	Introduction	59
6.2	Model procedures manuals and instruments.....	59
6.2.1	Background questions	60
6.2.2	Tasks Items	60
6.2.3	Standardized non-response coding.....	60

Table of Contents (cont)

6.3	Scoring.....	63
6.3.1	Intra-country rescoring.....	63
6.3.2	Inter-country rescoring.....	64
6.4	Data capture, data processing and coding.....	66
6.5	Derived Variables	66
7.0	Guidelines for Tabulation and Analysis	67
7.1	Sample Weighting Guidelines for Tabulation	67
7.2	Definitions of Types of Estimates: Categorical vs. Quantitative	67
7.2.1	Tabulation of categorical estimates.....	68
7.2.2	Tabulation of Quantitative Estimates.....	68
7.3	Skill Level Estimates.....	69
7.4	Rounding Guidelines	70
8.0	Data Quality.....	72
8.1	Sampling Errors	72
8.1.1	CV Release Guidelines.....	73
8.1.2	Using Plausible Values and Replicate Weights to calculating Sampling Error	75
8.1.3	Estimating Error Variance in ALL	83
8.1.4	Performing Analyses with the ALL Data Using SPSS	90
8.1.5	Performing Analyses with the ALL Data Using SAS	106
8.2	Non-Sampling errors.....	112
8.2.1	Sampling Frame	112
8.2.2	Non-response	112
8.2.3	Response Error.....	112
8.2.4	Scoring.....	113
9.0	Record Layouts and Univariate Counts	114
10.0	Principal Participants in the Project.....	115

1.0 Introduction

The data file in this package is a compilation of the ALL datasets received for a group of seven countries or regions that collected data in 2003. They include Bermuda, Canada, Italy, Norway, Switzerland, the United States and the Mexican State of Nuevo Leon. Another group of three countries collected their data in 2006 or 2008. They include Hungary, Netherlands, New Zealand. This document summarizes the survey concepts and operations of the international survey. It is important for users to become familiar with the contents of this document before publishing or otherwise releasing any estimates derived from the ALL microdata file.

For questions concerning the ALL microdata file, please contact:

CTCES Client Services

Client Services, Centre for Education Statistics
Statistics Canada,
150 Tunney's Pasture Driveway
Ottawa ON K1A 0T6
Tel: 1-800-307-3382
Fax : 613-951-4441
E-mail: educationstats@statcan.ca

2.0 Survey Overview

The Adult Literacy and Life Skills Survey (ALL) is a large-scale co-operative effort undertaken by governments, national statistics agencies, research institutions and multi-lateral agencies. The development and management of the study were co-ordinated by Statistics Canada and the Educational Testing Service (ETS) in collaboration with the National Center for Education Statistics (NCES) of the United States Department of Education, the Organisation for Economic Cooperation and Development (OECD), the Regional Office for Latin America and the Caribbean (OREALC) and the Institute for Statistics (UIS) of the United Nations Educational, Scientific and Cultural Organisation (UNESCO).

The survey instruments were developed by international teams of experts with financing provided by the Governments of Canada and the United States. A highly diverse group of countries and experts drawn from around the world participated in the validation of the instruments. Participating governments absorbed the costs of national data collection and a share of the international overheads associated with implementation.

The ALL study builds on the International Adult Literacy Survey (IALS), the world's first internationally comparative survey of adult skills undertaken in three rounds of data collection between 1994 and 1998. The foundation skills measured in the ALL survey include prose literacy, document literacy, numeracy, and problem solving. Additional skills assessed indirectly include familiarity with and use of information and communication technologies.

3.0 Survey Objectives

The ALL was initiated with two fundamental goals:

- 1) The first objective was to build on the skills measured in the 1994 IALS by introducing new assessment domains with robust theoretical frameworks and stable item parameters across countries and languages. This goal also involved directly linking the IALS with the ALL along the two literacy domains in order to allow comparison between prose and document profiles as measured in 1994 and later in 2003.
- 2) The second objective was to allow for, international, national and sub-national analysis of the correlates and possible antecedents of skills by collection a large enough bank of information from a sufficiently large number of respondents.

The central element of the survey was the direct assessment of the literacy, numeracy and problem solving skills of respondents using commonplace tasks of varying degree of difficulty drawn from a range of topic and knowledge areas. This information was supported by the collection of background information on respondents. In addition, the background questionnaire included questions on the self-assessment of literacy and numeracy skills of respondents, on the training which the respondent has taken in the year previous to the survey and on the perceived barriers to realizing enhanced literacy or numeracy skill levels.

4.0 Concepts and Definitions

4.1 Defining and Measuring Proficiency: Overview

For ALL, each proficiency scale starts at zero and increases to a theoretical maximum of 500 points. Scores along the scale denote the points at which a person with a given level of performance has an 80 percent probability of successfully completing a task at that level of difficulty. For instance, a person with an assessed performance at 250 points has an 80 percent probability of correctly answering a task with an estimated difficulty level of 250. The same individual would have an “80 percent plus” probability of correctly answering a simpler task (about 95% for a task with a complexity of 200) and a diminished probability (less than 80%) of successfully completing a more difficult task (about 40% for a task with a complexity of 300).¹

In an effort to facilitate analysis, these continuous scores have been regrouped into 5 skill levels (only 4 levels were defined for the problem solving scale) with level 1 being the lowest measured level of proficient. The proficiency levels used for ALL are useful in summarizing the results but also have some limitations. First, the relatively small proportions of respondents who actually reach Level 5 do not always allow for accurate reporting. For this reason, whenever results are presented by proficiency level, Levels 4 and 5 are typically combined. Second, the levels indicate specific sets of abilities and, therefore, the thresholds for the levels are not equidistant. The ranges of scores in each level are therefore not identical. In fact, for all four domains, Level 1 captures almost half of the scale. The thresholds for the problem solving domain are set somewhat differently and Level 1 covers precisely half of the scale. Level 1 includes all basic abilities required to attain higher levels. In other words, the ability to read may lie somewhere in Level 1, but the ability to understand and use what has been read comes in gradations of complexity from Level 1 to Level 5. The upshot of the relatively large ranges of scores in Level 1 on each of the scales is that there are multiple sub-levels of proficiency within this level. The range includes those who can barely read at all as well as those who read poorly or inattentively.²

Literacy and illiteracy

Interestingly, while the probability of a correct response may approach zero as the tasks become more difficult, it can never quite reach it because there is always some chance, however small, that a correct answer will be provided regardless of ability.

Accordingly, the results from the ALL measure performance along a proficiency continuum. The scales do not measure the absence of a competence, and thus cannot distinguish those who have from those who lack a specific competency.

This chapter offers a brief overview of the frameworks that were used to develop and interpret the scales used to measure prose and document literacy, numeracy, and problem solving in the International Adult Literacy and Skills Survey (ALL). The importance of developing a framework is thought to be central in construct-based approaches to measurement. Among the things that should be included in any such framework are an agreed upon definition of what ought to be measured and the identification of characteristics that can be used in the construction and interpretation of tasks. In addition to describing these characteristics for each measure, this

¹ Kirsch, Jungeblut and Campbell (1992), pp. 14-15.

² The International Survey of Reading Skills is a follow-up to the 2003 ALL that will provide more information about respondents at Level 1. Results are expected sometime in 2006.

chapter also includes sample items along with the identification of item features that are shown to contribute to item difficulty. Collectively this information provides a means for moving away from interpreting survey results in terms of discrete tasks or a single number and towards identifying levels of performance sufficiently generalized to have validity across assessments and groups.

4.2 Understand what was measured in the International Literacy and Skills Survey

4.2.1 Introduction

In 1992, the Organization for Economic Co-operation and Development (OECD) (OECD, 1992) concluded that low literacy levels were a serious threat to economic performance and social cohesion on an international level. But a broader understanding of literacy problems across industrialized nations – and consequent lessons for policy makers – was hindered due to a lack of comparable international data. Statistics Canada and Educational Testing Service (ETS) teamed up to build and deliver an international comparative study of literacy.

The International Adult Literacy Survey (IALS) was the first comparative survey of adults designed to profile and explore comparative literacy distributions among participating countries. In 2000, a final report was released (OECD and Statistics Canada, 2000) which included the results from three rounds of assessments involving some 23 country/language groups representing just over 50 per cent of the world's GDP. While IALS laid an important foundation for international comparative surveys of adults, there were also calls to expand what was being measured. There was a growing concern among governments and policy makers as to what additional competencies are relevant for an individual to participate fully and successfully in a modern society and for a society to meet the challenges of a rapidly changing world. One project aimed at addressing this issue was entitled *Definition and Selection of Key Competencies* (DeSeCo) and was carried out under the leadership of Switzerland. Its goal was to lay out, from a theoretical perspective, a set of key competencies that are believed to contribute to a successful life and a well-functioning society (Rychen and Salganik, 2003).

In response to these calls for broader measures, the ALL survey commissioned the development of frameworks to use as the basis for introducing new measures into the comparative assessments of adults. Those responsible for the development of ALL recognized that the design of any reliable and valid instrument should begin with a strong theoretical underpinning that is represented by a framework that characterizes current thinking in the field. According to Messick (1994) any framework that takes a construct-centered approach to assessment design should: begin with a general definition or statement of purpose – one that guides the rationale for the survey and what should be measured in terms of knowledge, skills or other attributes; identify various performances or behaviours that will reveal those constructs, and; identify task characteristics and indicate how these characteristics will be used in constructing the tasks that will elicit those behaviours.

4.2.2 Scaling the literacy, numeracy and problem solving tasks in ALL

The results of the ALL survey can be reported along four scales – two literacy scales (prose and document), a single numeracy scale, and a scale capturing problem solving – with each ranging from 0 to 500 points. One might imagine these tasks arranged along their respective scale in terms of their difficulty for adults and the level of proficiency needed to respond correctly to each task. The procedure used in ALL to model these continua of difficulty and ability is Item Response Theory (IRT). IRT is a mathematical model used for estimating the probability that a particular person will respond correctly to a given task from a specified pool of tasks (Murray, Kirsch and Jenkins, 1998).

The scale value assigned to each item results from how representative samples of adults in participating countries perform on each item and is based on the theory that someone at a given point on the scale is equally proficient in all tasks at that point on the scale. For the ALL survey, as for the IALS, proficiency was determined to mean that someone at a particular point on the proficiency scale would have an 80 per cent chance of answering items at that point correctly.

Just as adults within each participating country in ALL are sampled from the population of adults living in households, each task that was constructed and used in the assessment represents a type of task sampled from the domain or construct defined here. Hence, it is representative of a particular type of literacy, numeracy or problem solving task that is associated with adult contexts.

One obvious question that arises once one looks at the distributions of tasks along each of the described scales is, what distinguishes tasks at the lower end of each scale from those in the middle and upper ranges of the scale? Do tasks, that fall around the same place on each scale share some set of characteristics that result in their having similar levels of difficulty? Even a cursory review of the items reveals that tasks at the lower end of each scale differ from those at the higher end.

In an attempt to display this progression of complexity and difficulty, each proficiency scale was divided into levels. Both the literacy and numeracy scales used five levels where Level 1 represents the lowest level of proficiency and Level 5 the highest. These levels are defined as follows: Level 1 (0-225), Level 2 (226-275), Level 3 (276-325), Level 4 (326-375) and Level 5 (376-500). The scale for problem solving used four levels where Level 1 is the lowest level of proficiency and Level 4 the highest. These four levels are defined as follows: Level 1 (0-250), Level 2 (251-300), Level 3 (301-350), and Level 4 (351-500).

Since each level represents a progression of knowledge and skills, individuals within a particular level not only demonstrate the knowledge and skills associated with that level but the proficiencies associated with the lower levels as well. In practical terms, this means that individuals performing at 250 (the middle of Level 2 on one of the literacy or numeracy scales) are expected to be able to perform the average Level 1 and Level 2 task with a high degree of proficiency. A comparable point on the problem solving scale would be 275. In ALL, as in IALS, a high degree of proficiency is defined in terms of a response probability of 80 (RP80). This means that individuals estimated to have a particular scale score are expected to perform tasks at that point on the scale correctly with an 80 per cent probability. It also means they will have a

greater than 80 per cent chance of performing tasks that are lower on the scale. It does not mean, however, that individuals with given proficiencies can never succeed at tasks with higher difficulty values; they may do so some of the time. It does suggest that their probability of success is “relatively” low – i.e., the more difficult the task relative to their proficiency, the lower the likelihood of a correct response.

An analogy might help clarify this point. The relationship between task difficulty and individual proficiency is much like the high jump event in track and field, in which an athlete tries to jump over a bar that is placed at increasing heights. Each high jumper has a height at which he or she is proficient – that is, the jumper can clear the bar at that height with a high probability of success, and can clear the bar at lower heights almost every time. When the bar is higher than the athlete’s level of proficiency, however, it is expected that the athlete will be unable to clear the bar consistently.

4.3 Measuring prose and document literacy in ALL

The National Adult Literacy Survey (NALS), which was funded by the National Center for Education Statistics (NCES) as part of its overall assessment program in adult literacy, was the largest and most comprehensive study of adult literacy ever conducted in the United States (Kirsch, Jungeblut, Jenkins and Kolstad, 1993). Like all large-scale assessments funded by NCES, NALS was guided by a committee, which was comprised of a group of nationally recognized scholars, practitioners, and administrators who adopted the following definition of literacy:

“Literacy is using printed and written information to function in society, to achieve one’s goals, and to develop one’s knowledge and potential.”

This definition captures the initial work of the committee guiding the development of the assessment and provides the basis for creating other aspects of the framework to be discussed. It was also reviewed and adopted by the countries participating in the first round of IALS and was carried forward in ALL. This definition includes several assumptions made by panel members and, thus, it is important to consider various parts of this definition in turn.

Beginning with “**Literacy is...**”, the term literacy is used in preference to “reading” because it is likely to convey more precisely to a non-expert audience what the survey is measuring. “Reading” is often understood as simply decoding, or reading aloud, whereas the intention of the adult surveys is to measure something broader and deeper. Researchers studying literacy within particular contexts noted that different cultures and groups may value different kinds of literacy practices (Sticht, 1975; Heath, 1980; Szwed, 1981). Heath, for example, found that uses for reading could be described in terms of instrumental, social interactional, news-related, memory supportive, substitutes for oral messages, provision of a permanent record, and personal confirmation. The fact that people read different materials for different purposes implies a range of proficiencies that may not be well captured by signing one’s name, completing a certain number of years of schooling, or scoring at an 8th-grade level on a test of academic reading comprehension.

The phrase “... using printed and written information” draws attention to the fact that panel members view literacy not as a set of isolated skills associated with reading and writing, but more importantly as the application of those skills for specific purposes in specific contexts. When literacy is studied within varying contexts, diversity becomes its hallmark. First, people engage in literacy behaviours for a variety of uses or purposes (Sticht, 1978; Heath, 1980; Cook-Gumperz and Gumperz, 1981; Mikulecky, 1982). These uses vary across contexts (Heath, 1980; Venezky, 1983) and among people within the same context (Kirsch and Guthrie, 1984a). This variation in use leads to an interaction with a broad range of materials that have qualitatively different linguistic forms (Diehl, 1980; Jacob, 1982; Miller, 1982). In some cases, these different types of literacy tasks have been associated with different cognitive strategies or reading behaviours (Sticht, 1978, 1982; Crandall, 1981; Scribner and Cole, 1981; Kirsch and Guthrie, 1984b).

The phrase “... to function in society, to achieve one’s goals, and to develop one’s knowledge and potential ” is meant to capture the full scope of situations in which literacy plays a role in the lives of adults, from private to public, from school to work, to lifelong learning and active citizenship. “To achieve one’s goals and to develop one’s knowledge and potential” points to the view that literacy enables the fulfillment of individual aspirations—those that are defined such as graduation or obtaining a job, and those less defined and less immediate which extend and enrich one’s personal life. The phrase “to function in society” is meant to acknowledge that literacy provides individuals with a means of contributing to as well as benefiting from society. Literacy skills are generally recognized as important for nations to maintain or improve their standard of living and to compete in an increasingly global market place. Yet, they are equally as important for individual participation in technologically advancing societies with their formal institutions, complex legal systems, and large government programs.

4.3.1 Identifying task characteristics

The task characteristics represent variables that can be used in a variety of ways in developing an assessment and interpreting the results. Almond and Mislevy (1998) have identified five roles that variables can take on. They can be used to limit the scope of the assessment, characterize the features that should be used for constructing tasks, control the assembly of tasks into booklets or test forms, characterise examinees’ performance on or responses to tasks, or help to characterise aspects of competencies or proficiencies. IALS focused on variables that can be used to help in the construction of tasks as well as in the characterization of performance along one or more proficiency scales.

Each task in the assessment represents a piece of evidence about a person’s literacy (Mislevy, 2000). While the goal of the assessment will be to develop the best possible picture of an individual’s skills and abilities, the test cannot include an infinite number of tasks nor can an infinite number of features of those tasks be manipulated. Therefore, decisions need to be made about which features should be part of the test development process. Three task characteristics were identified and used in the construction of tasks for the IALS. These characteristics include:

Adult contexts/content. Since adults do not read written or printed materials in a vacuum, but read within a particular context or for a particular purpose, materials for the literacy assessment are selected that represent a variety of contexts and contents. This is to help ensure that no one

group of adults is either advantaged or disadvantaged due to the context or content included in the assessment. Six adult context/content categories have been identified as follows:

- ❖ **Home and family:** may include materials dealing with interpersonal relationships, personal finance, housing, and insurance.
- ❖ **Health and safety:** may include materials dealing with drugs and alcohol, disease prevention and treatment, safety and accident prevention, first aid, emergencies, and staying healthy.
- ❖ **Community and citizenship:** may include materials dealing with staying informed and community resources.
- ❖ **Consumer economics:** may include materials dealing with credit and banking, savings, advertising, making purchases, and maintaining personal possessions.
- ❖ **Work:** may include materials that deal in general with various occupations but not job specific texts, finding employment, finance, and being on the job.
- ❖ **Leisure and recreation:** may include materials involving travel, recreational activities, and restaurants.

Materials/texts. While no one would doubt that a literacy assessment should include a range of material, what is critical to the design and interpretation of the scores that are produced are the range and specific features of the text material which are included in constructing the tasks. A key distinction among texts that is at the heart of the IALS survey is their classification into continuous and non-continuous texts. Conventionally, continuous texts are formed of sentences organized into paragraphs. In these texts, organization occurs by paragraph setting, indentation, and the breakdown of text into a hierarchy signalled by headings that help the reader to recognize the organization of the text. The primary classification of continuous texts is by rhetorical purpose or text type. For IALS, these included: expository, descriptive, argumentative, and injunctive.

Non-continuous texts are organized differently than continuous texts and so allow the reader to employ different strategies for entering and extracting information from them. On the surface, these texts appear to have many different organizational patterns or formats, ranging from tables and schedules to charts and graphs, and from maps to forms. However, the organizational pattern for these types of texts, which Mosenthal and Kirsch (1998) refer to as documents, is said to have one of four basic structures: a simple list; a combined list; an intersected list; and a nested list. Together, these four types of documents make up what they have called matrix documents, or non-continuous texts with clearly defined rows and columns. They are also closely related to other non-continuous texts that these authors refer to as graphic, locative, and entry documents.

The distinction between continuous and non-continuous texts formed the basis for two of the three literacy scales used in IALS. Continuous texts were the basis for tasks that were placed along the prose scale while non-continuous texts formed the basis for tasks along the document scale. The quantitative scale included texts that were both continuous and non-continuous. The distinguishing characteristic for this scale was that respondents needed to identify and perform one or more arithmetic operations based on information contained in the texts. This scale was replaced in ALL with the numeracy scale, which is discussed in more detail later in this annex.

Processes/strategies. This task characteristic refers to the way in which examinees process text to respond correctly to a question or directive. It includes the processes used to relate information in the question (the given information) to the necessary information in the text (the new information) as well as the processes needed to either identify or construct the correct response from the information available. Three variables used to investigate tasks from national and international surveys will be summarized here. These are: type of match, type of information requested, and plausibility of distracting information.

4.3.2 Type of match

Four types of matching strategies were identified: locating, cycling, integrating, and generating. **Locating** tasks require examinees to match one or more features of information stated in the question to either identical or synonymous information provided in the text. **Cycling** tasks also require examinees to match one or more features of information, but unlike locating tasks, they require respondents to engage in a series of feature matches to satisfy conditions stated in the question.

Integrating tasks require examinees to pull together two or more pieces of information from the text according to some type of specified relation. For example, this relation might call for examinees to identify similarities (i.e., make a comparison), differences (i.e., contrast), degree (i.e., smaller or larger), or cause-and-effect relations. This information may be located within a single paragraph or it may appear in different paragraphs or sections of the text. In integrating information, examinees draw upon information categories provided in a question to locate the corresponding information in the text. They then relate the text information associated with these different categories based upon the relation term specified in the question. In some cases, however, examinees must *generate* these categories and/or relations before integrating the information stated in the text.

In addition to requiring examinees to apply one of these four strategies, the type of match between a question and the text is influenced by several other processing conditions which contribute to a task's overall difficulty. The first of these is the number of phrases that must be used in the search. Task difficulty increases with the amount of information in the question for which the examinee must search in the text. For instance, questions that consist of only one independent clause tend to be easier, on average, than those that contain several independent or dependent clauses. Difficulty also increases with the number of responses that examinees are asked to provide. Questions that request a single answer are easier than those that require three or more answers. Further, questions which specify the number of responses tend to be easier than those that do not. For example, a question which states, "List the 3 reasons..." would be easier than one which said, "List the reasons...". Tasks are also influenced by the degree to which examinees have to make inferences to match the given information in a question to corresponding information in the text, and to identify the requested information.

4.3.3 Type of information requested

This refers to the kinds of information that readers need to identify to answer a test question successfully. The more concrete the requested information, the easier the task is judged to be.

In previous research based on large-scale assessments of adults' and children's literacy (Kirsch and Mosenthal, 1994; Kirsch, Jungeblut, and Mosenthal, 1998), the type of information variable was scored on a 5-point scale. A score of one represented information that was the most concrete and therefore the easiest to process, while a score of five represented information that was the most abstract and therefore the most difficult to process.

For instance, questions which asked examinees to identify a person, animal, or thing (i.e., imaginable nouns) were said to request highly concrete information and were assigned a value of one. Questions asking respondents to identify goals, conditions, or purposes were said to request more abstract types of information. Such tasks were judged to be more difficult and received a value of three. Questions that required examinees to identify an "equivalent" were judged to be the most abstract and were assigned a value of five. In such cases, the equivalent tended to be an unfamiliar term or phrase for which respondents had to infer a definition or interpretation from the text.

4.3.4 Plausibility of distractors

This concerns the extent to which information in the text shares one or more features with the information requested in the question but does not fully satisfy what has been requested. Tasks are judged to be easiest when no distractor information is present in the text. They tend to become more difficult as the number of distractors increases, as the distractors share more features with the correct response, and as the distractors appear in closer proximity to the correct response. For instance, tasks tend to be judged more difficult when one or more distractors meet some but not all of the conditions specified in the question and appear in a paragraph or section of text other than the one containing the correct answer. Tasks are judged to be most difficult when two or more distractors share most of the features with the correct response and appear in the same paragraph or node of information as the correct response.

4.3.5 Characterizing prose literacy tasks

There are 55 tasks ordered along the 500-point prose literacy scale representing 19 IALS prose literacy tasks and 36 new prose literacy tasks designed and developed for the ALL survey. These tasks range in difficulty value from 169 to 439. One of the easiest tasks (receiving a difficulty value of 188 and falling in Level 1) directs the reader to look at a medicine label to determine the "maximum number of days you should take this medicine." Predictably, this item was also used as one of the contributing stimuli for the Health Literacy domain. In terms of our process variables, type of match was scored as easy because the reader was required to locate a single piece of information that was literally stated in the medicine label. The label contained only one reference to number of days and this information was located under the label dosage. Type of information was scored as easy because it asked for a number of days and plausibility of distractor was judged to be easy because there is no other reference to days in the medicine label.

MEDCO ASPIRIN

500

INDICATIONS: Headaches, muscle pains, rheumatic pains, toothaches, earaches. RELIEVES COMMON COLD SYMPTOMS.

DOSAGE: ORAL. 1 or 2 tablets every 6 hours, preferably accompanied by food, for not longer than 7 days. Store in a cool, dry place.

CAUTION: Do not use for gastritis or peptic ulcer. Do not use if taking anticoagulant drugs. Do not use for serious liver illness or bronchial asthma. If taken in large doses and for an extended period, may cause harm to kidneys. Before using this medication for chicken pox or influenza in children, consult with a doctor about Reyes Syndrome, a rare but serious illness. During lactation and pregnancy, consult with a doctor before using this product, especially in the last trimester of pregnancy. If symptoms persist, or in case of an accidental overdose, consult a doctor. Keep out of reach of children.

INGREDIENTS: Each tablet contains
500 mg acetylsalicylic acid.
Excipient c.b.p. 1 tablet.
Reg. No. 88246



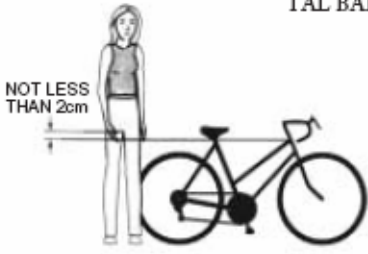
Made in Canada by STERLING PRODUCTS, INC.
1600 Industrial Blvd., Montreal, Quebec H9J 3P1

* Reprinted by permission

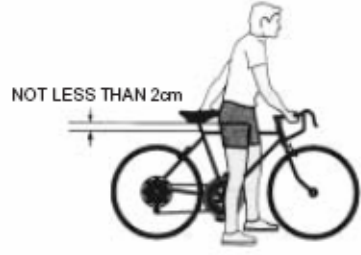
A second prose literacy task directs the reader to look at an article about impatiens. This task FALL in the middle of Level 2 and has a difficulty value of 254. It asks the reader to identify “what the smooth leaf surfaces and the stems suggest about the plant.” Again, the task directed the reader to locate information contained in the text so it was scored easy for type of information. The last sentence in the second paragraph under the heading **Appearance** states: “The smooth leaf surfaces and the stems indicate a great need of water.” Type of information was scored as being moderate because it directs the reader to identify a condition. Plausibility of distractor was scored as being moderate also because the same paragraph contained a sentence which serves to distract a number of readers. This sentence states, “... stems are branched and very juicy, which means, because of the tropical origin, that the plant is sensitive to cold.”

PROPER FRAME FIT

RIDER MUST BE ABLE TO STRADDLE BICYCLE WITH AT LEAST 2 cm CLEARANCE ABOVE THE HORIZONTAL BAR WHEN STANDING.



NOT LESS THAN 2cm



NOT LESS THAN 2cm

NOTE: Measurement for a female should be determined using a men's model as a basis.

PROPER SIZE OF BICYCLE

FRAME SIZE	LEG LENGTH OF RIDER
430mm	660mm-760mm
460mm	690mm-790mm
480mm	710mm-790mm
530mm	760mm-840mm
560mm	790mm-860mm
580mm	810mm-890mm
635mm	860mm-940mm

OWNER'S RESPONSIBILITY

1. **Bicycle Selection and Purchase:** Make sure this bicycle fits the intended rider. Bicycles come in a variety of sizes. Personal adjustment of seat and handlebars is necessary to assure maximum safety and comfort. Bicycles come with a wide variety of equipment and accessories . . . make sure the rider can operate them.
2. **Assembly:** Carefully follow all assembly instructions. Make sure that all nuts, bolts and screws are securely tightened.
3. **Fitting the Bicycle:** To ride safely and comfortably, the bicycle must fit the rider. Check the seat position, adjusting it up or down so that with the sole of rider's foot on the pedal in its lowest position the rider's knee is slightly bent. Note: Specific charts illustrated at left detail the proper method of determining the correct frame size.

The manufacturer is not responsible for failure, injury, or damage caused by improper completion of assembly or improper maintenance after shipment.

Tasks which fall at higher levels along the scale present the reader with more varied demands in terms of the type of match that is required and in terms of the number and nature of distractors that are present in the text. One such task (with a difficulty value of 281 or the beginning of Level 3) refers the reader to a page from a bicycle's owner's manual to determine how to ensure the seat is in the proper position. Type of information was scored as moderate because the reader needed to identify and state two conditions that needed to be met in writing. In addition, they were not told how many features they needed to provide from among those stated. Type of information was also scored as moderate also because it involved identifying a condition and plausibility of distractor received a score indicating it was relatively easy.

A somewhat more difficult task (318), one near the top of Level 3, involves an article about cotton diapers and directs the reader to "list three reasons why the author prefers to use disposable rather than cotton diapers." This task is made more difficult because of several of our process variables. First, type of match was scored as difficult because the reader had to provide multiple responses, each of which required a text-based inference. Nowhere in the text does the author say, "I prefer cotton diapers because...". These inferences are made somewhat more difficult because the type of information being requested is a "reason" rather than

something more concrete. This variable also was coded as difficult because of its abstractness. Finally, plausibility of distractor was scored as moderate because the text contains information that may serve to distract the reader.

An additional task falling in Level 4 on the Prose literacy scale (338) directs the reader to use the information from a pamphlet about hiring interviews to “write in your own words one difference between the panel and the group interview.” Here the difficulty does not come from locating information in the text. Rather than merely locating a fact about each type of interview, the reader needs to integrate what they have read to infer a characteristic on which the two types of interviews differ. Experience from other surveys of this kind reveal that tasks in which readers are asked to contrast information are more difficult, on average, than tasks in which they are asked to find similarities. Thus, type of match was scored as complex and difficult. Type of information was scored as being difficult as well because it directs the reader to provide a difference. Differences tend to be more abstract in that they ask for the identification of distinctive or contrastive features related in this case to an interview process. Plausibility of distractor was judged as being easy because no distracting information was present in the text. Thus this variable was not seen as contributing to the overall difficulty of this task.

The Hiring Interview

Preinterview

Try to learn more about the business. What products does it manufacture or services does it provide? What methods or procedures does it use? This information can be found in trade directories, chamber of commerce or industrial directories, or at your local employment office.

Find out more about the position. Would you replace someone or is the position newly created? In which departments or shops would you work? Collective agreements describing various standardized positions and duties are available at most local employment offices. You can also contact the appropriate trade union.

The Interview

Ask questions about the position and the business. Answer clearly and accurately all questions put to you. Bring along a note pad as well as your work and training documents.

The Most Common Types of Interview

One-on-one: Self explanatory.

Panel: A number of people ask you questions and then compare notes on your application.

Group: After hearing a presentation with other applicants on the position and duties, you take part in a group discussion.

Postinterview

Note the key points discussed. Compare questions that caused you difficulty with those that allowed you to highlight your strong points. Such a review will help you prepare for future interviews. If you wish, you can talk about it with the placement officer or career counsellor at your local employment office.

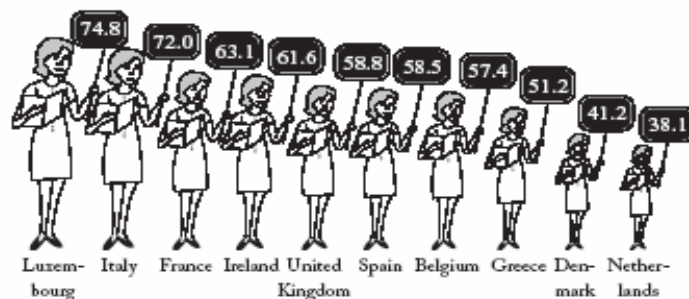
The most difficult task on the prose literacy scale (377) fALL in the lower range of Level 5 and required readers to look at an announcement from a personnel department and to “list two ways in which CIEM (an employee support initiative within a company) helps people who lose their jobs because of departmental reorganization.” Type of match was scored difficult because the question contained multiple phrases that the reader needed to keep in mind when reading the text. In addition, readers had to provide multiple responses and make low text-based inferences. Type of information received a moderate score because readers were looking for a purpose or function and plausibility of distractor was scored as relatively difficult. This task is made somewhat more difficult because the announcement is organized around information that is different from what is being requested in the question. Thus while the correct information is listed under a single heading, this information is embedded under a list of headings describing CIEM’s activities for employees looking for other work. Thus, this list of headings in the text serves as an excellent set of distractors for the reader who does not search for or locate the phrase in the question containing the conditional information – those who lose their jobs because of a departmental reorganization.

4.3.6 Characterizing document literacy tasks

There are 54 tasks ordered along the 500-point document literacy scale. These 54 tasks comprise 19 items from IALS and 35 new tasks developed for ALL. Together, these tasks range in difficulty value from 157 to 444. A Level 1 document literacy task with a difficulty value of 188 directs the reader to identify from a chart the percentage of teachers from Greece who are women. The chart shown here displays the percentage of teachers from various countries who are women. In terms of our process variables, type of match was judged to be easy because the reader was required to locate a single piece of information that was literally stated in the chart; type of information was judged to be relatively easy because it was an amount; and plausibility of distractor is also judged to be relatively easy because there are distractors for the requested information.

FEW DUTCH WOMEN AT THE BLACKBOARD

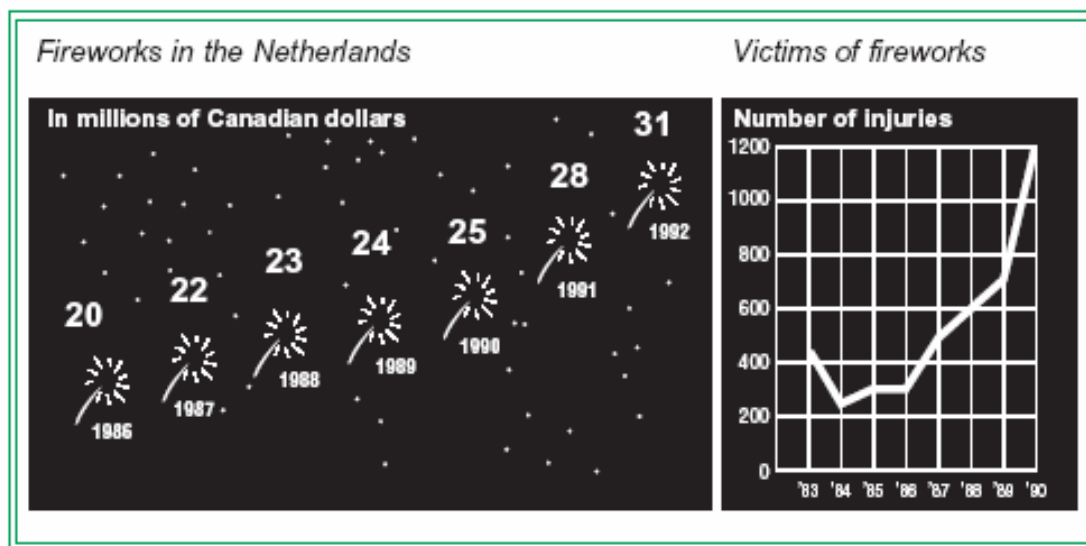
There is a low percentage of women teachers in the Netherlands compared to other countries. In most of the other countries, the majority of teachers are women. However, if we include the figures for inspectors and school principals, the proportion shrinks considerably and women are in a minority everywhere.



Percentage of women teachers (kindergarten, elementary, and secondary).

A second document task involving this same chart directs the reader to identify the country other than the Netherlands in which women teachers are in the minority. This item FALL in the middle of Level 2 and received a difficulty value of 234. This task was made a bit more difficult than the first because rather than searching for a country and locating a percentage, the reader had to know that minority means less than 50 per cent. Then they had to cycle through to identify the countries in which the percentage of women teachers were less than 50 per cent. In addition, they had to remember the condition “other than the Netherlands”; otherwise they might have chosen it over the correct response. As a result, type of match was scored as moderately difficult; type of information as easy because the requested information is a country or place; and plausibility of distractor as relatively easy because there are distractors associated with the requested information.

A somewhat more difficult task, with a difficulty value of 295 and falling in the middle of Level 3 directs the reader to look at charts involving fireworks from the Netherlands and to write a brief description of the relationship between sales and injuries based on the information shown. Here the reader needs to look at and compare the information contained in the two charts and integrate this information making an inference regarding the relationship between the two sets of information. As a result, it was judged as being relatively difficult in terms of type of match. Type of information also was judged to be relatively difficult because the requested information is asking for a pattern or similarity in the data. Plausibility of distractor was scored moderately difficult primarily because both given and requested information is present in the task. For example, one of the things that may have contributed to the difficulty of this task is the fact that the sales graph goes from 1986 to 1992 while the injuries graph goes from 1983 to 1990. The reader needed to compare the information from the two charts for the comparable period time.



Another set of tasks covering a range of difficulty on the document scale involved a rather complicated document taken from a page in a consumer magazine rating clock radios. The easiest of the three tasks, receiving a difficulty value of 287 and falling in Level 3, asks the reader “which two features are not on any basic clock radio.” In looking at the document, the reader has to cycle through the document, find the listing for basic clock radios, and then determine that a dash represents the absence of a feature. They then have to locate the two features indicated by the set of dashes. As a result, type of match was judged as being relatively difficult because it is a cycle requiring multiple responses with a condition or low text based inference. Type of information was scored as relatively easy because its features are an attribute of the clock radio and plausibility of distractor is relatively easy because there are some characteristics that are not associated with other clock radios.

A somewhat more difficult task associated with this document and falling in the lower end of Level 4 received a difficulty value of 327. It asks the reader “which full-featured clock radio is rated highest on performance.” Here the reader must make a three-feature match (full-featured, performance, and highest) where one of the features requires them to process conditional information. It is possible, for example, that some readers were able to find the full-featured radios and the column listed under performance but selected the first radio listed assuming it was the one rated highest. In this case, they did not understand the conditional information which is a legend stating what the symbols mean. Others may have gone to the column labelled overall score and found the highest numerical number and chosen the radio associated with it. For this reason, plausibility of distractor was scored as moderately difficult. Type of information was judged as being easy because the requested information is a thing.

The most difficult task associated with this document, with a difficulty level of 408, and falling in Level 5 asks the reader to identify the average advertised price for the basic clock radio receiving the highest overall score. This task was made more difficult because the reader had to match four rather than three features; they also had to process conditional information and there was a highly plausible distractor in the same node as the correct answer. As a result of these

18

RATINGS

4 Overall score: If A composite, encompassing all our tests and judgments. A 'perfect' radio would have earned 100 points.

2 Sensitivity. How well each radio received a station with little interference.

21 **Dual alarm.** Lets you set two separate wake-up times.

[1] Dacotefraud. Replaced by HC-X200, \$79 list and \$60 average advertised sale price.

L - Warranty repairs cost \$5 for handling.
M - Warranty repairs cost \$10 for handling.

4.4 Measuring numeracy in ALL

The conception of numeracy developed for ALL is built upon recent research and work done in several countries on functional demands of different life contexts, on the nature of adults' mathematical and statistical knowledge and skills, and on how such skills are applied or used in different circumstances. In light of the general intention of the ALL survey to provide information about a diverse set of life skills, this framework defines numeracy as follows:

Numeracy is the knowledge and skills required to effectively manage and respond to the mathematical demands of diverse situations.

This definition implies that numeracy is broader than the construct of quantitative literacy defined by IALS. Further, adult numeracy should be viewed as different from “knowing school mathematics”. Although a universally accepted definition of “numeracy” does not exist (Baker and Street, 1994), an examination of some perspectives on the meaning of adult numeracy shows that they contain many commonalities. Below are two examples, both from work in Australia:

Numeracy is the mathematics for effective functioning in one's group and community, and the capacity to use these skills to further one's own development and of one's community (Beazley, 1984).

Numeracy involves abilities that include interpreting, applying and communicating mathematical information in commonly encountered situations to enable full, critical and effective participation in a wide range of life roles (Queensland Department of Education, 1994)

All these definitions are quite similar, in their broad scope, to the ALL definitions of prose and document literacy presented in a prior section. Many conceptions of numeracy emphasize the practical or functional application and use of mathematical knowledge and skills to cope with the presence of mathematical elements in real situations. Adults are expected to possess multiple ways of responding flexibly to a mathematical situation in a goal-oriented way, dependent on the needs and interests of the individual within the given context (i.e., home, community, workplace, etc...), as well as on his or her attitudes and beliefs toward numeracy (Gal, 2000; Coben, O'Donoghue and FitzSimons, 2000).

Thus, numeracy involves more than just applying arithmetical skills to information embedded in printed materials, which was the focus of assessment in IALS. Adult numeracy extends to a possession of number sense, estimation skills, measurement and statistical literacy. Given the extent to which numeracy pervades the modern world, it is not necessarily just commonly encountered situations that require numerate behaviour, but also *new* situations.

Another important element in defining numeracy is the role of communication processes. Numeracy not only incorporates the individual's abilities to use and apply mathematical skills efficiently and critically, but also requires the person to be able to interpret textual or symbolic messages as well as to communicate mathematical information and reasoning processes (Marr and Tout, 1997; Gal, 1997).

Definitions of numeracy explicitly state that numeracy not only refers to operating with numbers, as the word can suggest, especially to those familiar with conceptions of children's numeracy, but covers a wide range of mathematical skills and understandings. Further, in recent years there has been much discussion and debate about the relationship between mathematics and numeracy and about the concept of "critical" numeracy (Frankenstein, 1989; Steen, 2001). Johnston, for example, has argued that:

To be numerate is more than being able to manipulate numbers, or even being able to 'succeed' in school or university mathematics. Numeracy is a critical awareness which builds bridges between mathematics and the real-world, with all its diversity (Johnston, 1994).

Many authors argue that a discussion of functional skills should also address supporting or enabling attitudes and beliefs. In the area of adults' mathematical skills, "at homeness" with numbers or "confidence" with mathematical skills is expected, as these affect how skills and knowledge are actually put into practice (Cockroft, 1982; Tobias, 1993).

The brief definition of numeracy developed for ALL and presented earlier above is complemented by a broader definition of **numerate behaviour** which was developed by the ALL Numeracy Team to serve as the basis for the development of numeracy items for ALL:

Numerate behaviour is observed when people manage a situation or solve a problem in a real context; it involves responding to information about mathematical ideas that may be represented in a range of ways; it requires the activation of a range of enabling knowledge, factors and processes.

This conception of numerate behaviour implies that in order to assess people's numeracy, it is necessary to generate tasks and items which vary in terms of contexts, the responses called for, the nature of the mathematical information involved, and the representations of this information. These task characteristics are elaborated below. This conception is much broader than the definition of quantitative literacy used in IALS. Its key elements relate in a broad way to situation management and to a need for a range of responses (not only to responses that involve numbers). It refers to a wide range of skills and knowledge (not only to application of arithmetical knowledge and computational operations) and to the use of a wide range of situations that present respondents with mathematical information of different types (not only those involving **numbers** embedded in **printed** materials).

The item development process aimed to ensure that a certain proportion of the item pool would place a minimum reading burden on the respondents, i.e., that some of the stimuli would be text-free or almost so, allowing even respondents with limited mastery of the language of the test to comprehend the situation described. Other parts of the item pool included items requiring varying amounts of essential texts as dictated by the situation which the item aimed to represent.

As implied by the literature and ideas reviewed earlier, the nature of a person's responses to the mathematical and other demands of a situation will depend critically on the activation of various

enabling knowledge bases (understanding of the context; knowledge and skills in the areas of mathematics, statistics and literacy), on reasoning processes and on their attitudes and beliefs with respect to numeracy. In addition, numerate behaviour requires the integration of mathematical knowledge and skills with broader literacy and problem solving skills along with the prior experiences and practices that each person brings to every situation. It is clear that numerate behaviour will involve an attempt to engage with a task and not delegate it to others or deal with it by intentionally ignoring its mathematical content.

4.4.1 Identifying task characteristics

Four key characteristics of numerate behaviour were used to develop and represent the numeracy tasks built for ALL – type of purpose/context, type of response, type of mathematical or statistical information, and type of representation of mathematical or statistical information. Each of these is described next.

Type of purpose/context. People try to manage or respond to a numeracy situation because they want to satisfy a purpose or reach a goal. Four types of purposes and goals are described below. To be sure, these are not mutually exclusive and may involve the same underlying mathematical themes.

4.4.2 Everyday life

The numeracy tasks that occur in everyday situations are often those that one faces in personal and family life, or revolve around hobbies, personal development, or interests. Representative tasks are handling money and budgets, comparison shopping, planning nutrition, personal time management, making decisions involving travel, planning trips, mathematics involved in hobbies like quilting or wood-working, playing games of chance, understanding sports scoring and statistics, reading maps and using measurements in home situations such as cooking or home repairs.

4.4.3 Work-related

At work, one is confronted with quantitative situations that often are more specialized than those seen in everyday life. In this context, people have to develop skills in managing situations that might be narrower in their application of mathematical themes. Representative tasks are completing purchase orders, totalling receipts, calculating change, managing schedules, using spreadsheets, organizing and packing different shaped goods, completing and interpreting control charts or quality graphs, making and recording measurements, reading blueprints, tracking expenditures, predicting costs and applying formulas.

4.4.4 Societal or community

Adults need to know about processes happening in the world around them, such as trends in crime, wages and employment, pollution, medical or environmental risks. They may have to take part in social or community events, or in political action. This requires that adults can read and interpret quantitative information presented in the media, including statistical messages and graphs. They may have to manage situations like organizing a fund-raiser, planning fiscal

aspects of a community program, or interpreting the results of a study about risks of the latest health fad.

4.4.5 Further learning

Numeracy skills enable a person to participate in further study, whether for academic purposes or as part of vocational training. In either case, it is important to be able to know some of the more formal aspects of mathematics that involve symbols, rules, and formulas and to understand some of the conventions used to apply mathematical rules and principles.

Type of responses. In different types of real-life situations, people may have to respond in one or more of the following ways. (The first virtually always occurs; others will depend on the interaction between situational demands and the goals, skills, dispositions, and prior learning of the person):

Identify or locate some mathematical information present in the task or situation confronting them that is relevant to their purpose or goal.

Act upon or react to the information in the situation. Bishop (1988), for example, proposed that there are six modes of mathematical actions that are common in all cultures: counting, locating, measuring, designing, playing and explaining. Other types of actions or reactions may occur, such as doing some calculations (“in the head” or with a calculator), ordering or sorting, estimating, measuring, or modeling (such as by using a formula).

Interpret the information embedded within the situation (and the results of any prior action) and comprehend what it means or implies. This can include making a judgment about how mathematical information or known facts actually apply to the situation or context. Contextual judgment may have to be used in deciding whether an answer makes sense or not in the given context, for example, that a result of “2.35 cars” is not a valid solution to how many cars are needed to transport a group. It can also incorporate a critical aspect, where a person questions the purpose of the task, the validity of the data or information presented, and the meaning and implications of the results, both for them as an individual and possibly for the wider community.

Communicate about the mathematical information given, or the results of one’s actions or interpretations to someone else. This can be done orally or in writing (ranging from a simple number or word to a detailed explanation or analysis) and/or through drawing (a diagram, map, graph).

Type of mathematical or statistical information. Mathematical information can be classified in a number of ways and on different levels of abstraction. One approach is to refer to fundamental “big ideas” in the mathematical world. Steen (1990), for example, identified six broad categories pertaining to: quantity, dimension, pattern, shape, uncertainty, and change. Rutherford and Ahlgren (1990) described networks of related ideas: numbers, shapes, uncertainty, summarizing data, sampling and reasoning. Dossey (1997) categorized the

mathematical behaviours of quantitative literacy as: data representation and interpretation, number and operation sense, measurement, variables and relations, geometric shapes and spatial visualization, and chance. The ALL Numeracy Team drew from these and other closely tied categorizations (e.g., National Council of Teachers of Mathematics, 2000) to arrive at a set of five fundamental ideas that characterize the mathematical demands facing adults in diverse situations at the beginning of the 21st century.

4.4.6 Quantity and number

Quantity is described by Fey (1990) as an outgrowth of people's need to quantify the world around us, using attributes such as: length, area and volume of rivers or land masses; temperature, humidity and pressure of our atmosphere; populations and growth rates of species; motions of tides; revenues or profits of companies, etc...

Number is fundamental to quantification and different types of number constrain quantification in various ways: whole numbers can serve as counters or estimators; fractions, decimals and per cents as expressions of greater precision, or as indications of parts-of-whole which are useful when comparing proportions. Positive and negative numbers serve as directional indicators. In addition to quantification, numbers are used to put things in order and as identifiers (e.g., telephone numbers or zip codes). Facility with quantity, number, and operation on number requires a good "sense" for magnitude and the meaning of very large or very small numbers, and sometimes a sense for the relative magnitude of different proportions.

Money and time management, the ubiquitous mathematics that is part of every adult's life, depends on a good sense of number and quantity. Contextual judgment comes into play when deciding how precise one should be when conducting certain computations or affects the choice of which tool (calculator, mental math, a computer) to use. A low level numeracy task might be figuring out the cost of one can of soup, given the cost of four for \$2.00; a task with a higher cognitive demand could involve "harder numbers" such as when figuring out the cost per kilo while buying 0.783 kg of cheese for 12,95 Euros.

4.4.7 Dimension and shape

Dimension includes "big ideas" related to one, two and three dimensions of "things". Understanding of dimensions is called for when encountering or generating spatial or numerical descriptions of objects, making projections, or working with lengths, perimeters, planes, surfaces, location, etc... Facility with each dimension requires a sense of "benchmark" measures, direct measurement, and estimations of measurements.

Shape is a category describing real or imaginary images and entities that can be visualized (e.g., houses and buildings, designs in art and craft, safety signs, packaging, knots, crystals, shadows and plants). Direction and location are fundamental qualities called upon when reading or sketching maps and diagrams. A basic numeracy task in this fundamental aspect could be shape identification whereas a more complex task might involve describing the change in the size or volume of an object when one dimension is changed, such as when choosing between different boxes for packaging certain objects.

4.4.8 Pattern, functions and relationships

It is frequently written that mathematics is the study of patterns and relationships. Pattern is seen as a wide-ranging concept that covers patterns encountered all around us, such as those in musical forms, nature, traffic patterns, etc... It is argued by Senechal (1990) that our ability to recognize, interpret and create patterns is the key to dealing with the world around us. The human capacity for identifying relationships and for thinking analytically underlies mathematical thinking. Algebra – beyond symbolic manipulation – provides a tool for representing relationships between amounts through the use of tables, graphs, symbols and words. The ability to generalize and to characterize functions, relationships between variables, is a crucial gateway to understanding even the most basic economic, political or social analyses. A relatively simple pattern-recognition task might require someone to describe the pattern in a sequence of given numbers or shapes, and in a functional context to understand the relationship between lists or variables (e.g., weight and volume of objects); having to develop a formula for an electronic spreadsheet would put a higher level of demand on the individual.

4.4.9 Data and chance

Data and chance encompass two related but separate topics. **Data** covers “big ideas” such as variability, sampling, error, or prediction, and related statistical topics such as data collection, data analysis, and common measures of center or spread, or the idea of a statistical inference. Modern society demands that adults are able to interpret (and at times even produce) frequency tables, basic charts and graphs, information about averages and medians, as well as identify questionable statistical claims (Gal, 2002).

Chance covers “big ideas” related to probability and relevant statistical concepts and tools. Few things in the world are 100 per cent certain; thus the ability to attach a number that represents the likelihood of an event (including risks or side-effects) is a valuable tool whether it has to do with the weather, the stock-market, or the decision to use a certain drug. In this category, a simple numeracy skill might be the interpretation of a simple pie chart or comprehension of a statement about an average; a more complex task would be to infer the likelihood of occurrence of an event based upon given information.

4.4.10 Change

This term describes the mathematics of how the world changes around us. Individual organisms grow, populations vary, prices fluctuate, objects traveling speed up and slow down. Change and rates of change help provide a narration of the world as time marches on. Additive, multiplicative or exponential patterns of change can characterize steady trends; periodic changes suggest cycles and irregular change patterns connect with chaos theory. Describing weight loss over time is a relatively simple task, while calculating compounded interest is a relatively complex task.

Type of representation of mathematical information. Mathematical information in an activity or a situation may be available or represented in many forms. It may appear as concrete objects to be counted (e.g., sheep, people, buildings, cars, etc...) or as pictures of such things. It may be conveyed through symbolic notation (e.g., numerals, letters, or operation signs). Sometimes,

mathematical information will be conveyed by formulas, which are a model of relationships between entities or variables.

Further, mathematical information may be encoded in visual displays such as *diagrams* or **charts**; **graphs**, and **tables** may be used to display aggregate statistical or quantitative information. Similarly, **maps** of real entities (e.g., of a city or a project plan) may contain numerical data but also information that can be quantified or mathematized.

Finally, a person may have to extract mathematical information from various types of texts, either in prose or in documents with specific formats (such as in tax forms). Two different kinds of text may be encountered in functional numeracy tasks. The first involves mathematical information represented in textual form, i.e., with words or phrases that carry mathematical meaning. Examples are the use of number words (e.g., “five” instead of “5”), basic mathematical terms (e.g., fraction, multiplication, per cent, average, proportion), or more complex phrases (e.g., “crime rate cut by half”) that require interpretation. The second involves cases where mathematical information is expressed in regular notations or symbols (e.g., numbers, plus or minus signs, symbols for units of measure, etc...), but is surrounded by text that despite its non-mathematical nature also has to be interpreted in order to provide additional information and context. An example is a bank deposit slip with some text and instructions in which numbers describing monetary amounts are embedded.

4.4.11 Characterizing numeracy tasks

A total of 40 numeracy tasks were selected and used in the ALL survey. These tasks range along the numeracy scale from 174 to 380 and their placement was determined by how well adults in participating countries responded to each task. Described below are sample tasks that reflect some of the conceptual facets of the numeracy construct and scale design principles described earlier, such as computations, spatial and proportional reasoning, measurement, and statistical knowledge.

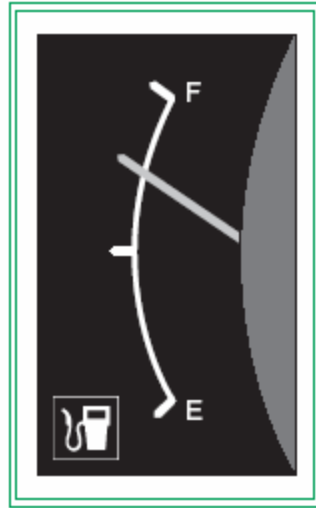
As expected, the easiest task on the numeracy scale required adults to look at a photograph containing two cartons of coca cola bottles (174). They were directed to find the total number of bottles in the two full cases being shown. Part of what made this task easy is the fact that content was drawn from everyday life and objects of this kind would be relatively familiar to most people. Second, what adults were asked to do was apparent and explicit – this task used a photograph depicting concrete objects and required the processing of no text. A third contributing factor is that respondents could approach the task in a variety of ways that differ in sophistication, such as by multiplying rows and columns, but also by simple counting. This task requires that adults make a conjecture since the full set of bottles in the lower case is not visible, but as can be seen from the low difficulty level of the task, this feature did not present a problem for the vast majority of adults in all participating countries.



A second task that was also quite easy directed adults to look at a short text depicting the results of an election involving three candidates and determine the total number of votes cast. This task received a difficulty value of 192, falling in Level 1 on the numeracy scale. Again, respondents were asked to deal with a realistic type of situation where simple numerical information is displayed in a simple column format showing the name of each candidate and the number of votes that the candidate received. No other numerical information was present that can be a distractor. Finding the total number of votes cast in the election requires a single addition operation that is made explicit in the question by the use of the keyword “total”, and the computation involves relatively small whole numbers.

A more complex numeracy task falling in the middle of Level 2 and receiving a difficulty value of 248 directs adults to look at a gas (petrol) gauge. This gauge has three lines or ticks on it with one showing an “F”, one showing an “E” and the third in the middle between the two. A line on the gauge, representing the gauge’s needle, shows a level that is roughly halfway between the middle tick and the tick indicating “F”, suggesting that the tank is about three-quarters full. The directive states that the tank holds 48 gallons and asks the respondent to determine “how many gallons remain in the tank.” This task is drawn from an everyday context and requires an adult to interpret a display that conveys quantitative information but carries virtually no text or numbers. No mathematical information is present other than what is given in the question.

What makes this task more difficult than the previous ones described above is the fact that adults must first estimate the level of gas remaining in the tank, by converting the placement of the needle to a fraction. Then they need to determine how many gallons this represents from the 48 gallon capacity stated in the question or directive. Thus, this task requires adults to apply multiple operations or procedures to arrive at a correct response, without specifying what the operations may be. Nonetheless, this task, like many everyday numeracy tasks, does not require an exact computation but allows an approximation that should fall within reasonable boundaries.



A somewhat more difficult numeracy task, falling at the top of Level 2 and receiving a difficulty value of 275, requires adults to look at a diagram of a container on which there are four markings or lines; respondents are asked to draw a line on the container indicating where one-third would be. The top line is marked “1” while the middle line is marked with “ $\frac{1}{2}$ ”. There are two other lines with no markings - one line midway between “1” and “ $\frac{1}{2}$ ” and another midway between the line marked “ $\frac{1}{2}$ ” and the bottom of the container. To respond correctly, adults need to mark a line on the container that is between the line marked “ $\frac{1}{2}$ ” and the line below it indicating where one-quarter would be (although this line does not say “ $\frac{1}{4}$ ” – this has to be inferred). Here the context may be less familiar to the respondent but again the visual image used is simple and realistic with virtually no text; the response expected does not involve writing a symbol or text, just drawing a line in a certain region on the drawing of the container. To answer this task correctly, adults need to have some working knowledge of fractions and a sense for proportions: they have to be familiar with the symbols for “ $\frac{1}{2}$ ” and “ $\frac{1}{3}$ ”, know how to order fractions in terms of their relative size and be able to relate them to the existing markings on the container.

Some numeracy tasks were developed around a short newspaper article titled “Is breast milk safe?” which relates to environmental hazards and food safety. The article contained two brief text paragraphs describing a toxin, Dioxin, found in fish in the Baltic Sea plus a graph with bars indicating the levels of Dioxin found at three points in time, namely 1975, 1985, and 1995, in the breast milk of North European women. One question asked adults to describe how the amount of Dioxin changed from 1975 to 1995, i.e., provide a straightforward interpretation of data presented in a graph. Adults were not required to actually calculate the amount of change over each of the periods, just describe in their own words the change in the levels of Dioxin (e.g., decreased, increased, stayed the same).

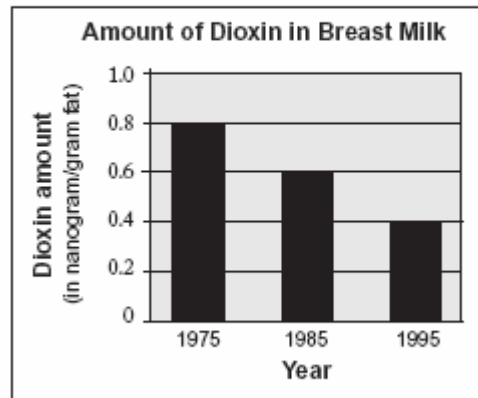
This task received a difficulty value of 280, the lower end of Level 3. The graph clearly indicates that the amount of Dioxin decreased over each of the three time periods, yet some adults have difficulty coping with such a task, which is based on a stimulus with a structure that commonly appears in newspapers, i.e., brief text plus a graph. The increased difficulty level of this item may be attributable in part to the need for adults to generate their own description, to the

moderate amount of dependence on text needed to comprehend the context to which the graph refers, or to the need to understand the direction of the decimal values on the vertical axis (which is common in reporting on concentrations of contaminating chemicals). This item also served to help create the Health Literacy Domain.

Is breast milk safe?

Since the 1970s, scientists have been worried about the amount of Dioxin, a toxin in fish caught in the Baltic sea. Dioxin tends to accumulate in breast milk and can harm newborn babies.

The diagram shows the amount of Dioxin in the breast milk of North European women, as found in studies done from 1975 to 1995.



A second and more difficult task using this same stimulus directed adults to compare the per cent of change in Dioxin level from 1975 to 1985 to the per cent of change in Dioxin level from 1985 to 1995, determine which per cent of change is larger, and explain their answer. This task was considerably more difficult for adults in participating countries and received a difficulty value of 377 on the numeracy scale. Here the necessary information is embedded within the graph and requires a level of transformation and interpretation. To arrive at a correct response, adults have to look at the rate of change expressed in per cents, not just the absolute size of the change. Further, they have to work with per cents of entities smaller than one (i.e., the decimal values on the vertical axis) and realize that the base for the computation of per cent change shifts for each pair. It seems that the need to cope with such task features, use formal mathematical procedures, or deal with the abstract notion of rate of change, adds considerable difficulty to such tasks.

The most difficult numeracy task in this assessment, receiving a difficulty value of 380 (Level 5), presented adults with an advertisement claiming that it is possible for an investor to double an amount invested in seven years, based on a 10 per cent fixed interest rate each year. Adults were asked if it is possible to double \$1000 invested at this rate after seven years and had to support their answer with their calculations. A range of responses was accepted as correct as long as a reasonable justification was provided, with relevant computations. Respondents were free to perform the calculation any way they wanted, but could also use a “financial hint” which accompanied the advertisement and presented a formula for estimating the worth of an investment after any number of years. Those who used the formula had to enter information stated in the text into variables in the formula (principal, interest rate and time period) and then perform the needed computations and compare the result to the expected amount if \$1000 is doubled.

All respondents could use a hand-held calculator provided as part of the assessment. This task proved difficult because it involved per cents and the computation, whether with or without the formula, required the integration of several steps and several types of operations. Performing the computations without the formula required understanding of compound interest procedures. This task allowed adults to use a range of reasoning strategies, including informal or invented procedures. Yet, like the previous task involving the comparison of rates of change, it required the use of formal mathematical information and deeper understanding of non-routine computational procedures, all of which may not be familiar or accessible to many adults.

4.5 Measuring problem solving in ALL

Research on problem solving has a long tradition within both academic psychology and applied human resources research. A very general definition of problem solving that reflects how it is generally understood in the psychological literature (Hunt, 1994; Mayer, 1992; Mayer and Wittrock, 1996; Smith, 1991) is presented here:

Problem solving is goal-directed thinking and action in situations for which no routine solution procedure is available. The problem solver has a more or less well-defined goal, but does not immediately know how to reach it. The incongruence of goals and admissible operators constitutes a problem. The understanding of the problem situation and its step-by-step transformation, based on planning and reasoning, constitute the process of problem solving.

One major challenge while developing a framework for problem solving that is to be used in a survey such as ALL is how best to adapt the psychological literature to the constraints imposed by a large-scale international comparative study. In order to do this, a decision was made to focus on an essential subset of problem solving – analytical problem solving. Our notion of analytical problem solving is not to be confused with the intuitive everyday use of the term or with the clinical-psychological concept in which problem solving is associated with the resolution of social and emotional conflicts. Nevertheless, social context is also relevant for our definition of analytical problem solving, for example when problems have to be approached interactively and resolved through co-operation. Motivational factors such as interest in the topic and task-orientation also influence the problem-solving process. However, the quality of problem solving is primarily determined by the comprehension of the problem situation, the thinking processes used to approach the problem, and the appropriateness of the solution.

The **problem** itself can be characterized by different aspects:

- ❖ The **context** can reflect different domains, which may be of a theoretical or a practical nature, related to academic situations or to the real world. Within these domains, problems can be more or less authentic.
- ❖ The **scope** of a problem can range from working on limited, concrete parts of a task to planning and executing complex actions or evaluating multiple sequences of actions.
- ❖ The problem can have a well-defined or an ill-defined goal, it can have transparent (explicitly named) or non-transparent constraints, and involve few independent elements

or numerous interconnected ones. These features determine the **complexity** of the problem.

How familiar the context is to the target population, whether the problem involves concrete tasks or complex actions, how well the goal is defined, how transparent the constraints are, how many elements the problem solver has to take into account and how strongly they are interconnected – are all features that will determine the level of problem-solving competency required to solve a certain problem. The empirical difficulty, i.e., the probability of giving a correct solution, will depend on the relation between these problem features on the one hand, and the subjects' competency level on the other hand.

The **cognitive processes** that are activated in the course of problem solving are diverse and complex, and they are likely to be organized in a non-linear manner. Among these processes, the following five components may be identified:

1. Searching for information, and structuring and integrating it into a mental representation of the problem (“situational model”).
2. Reasoning, based on the situational model.
3. Planning actions and other solution steps.
4. Executing and evaluating solution steps.
5. Continuous processing of external information and feedback.

Baxter and Glaser (1997) present a similar list of cognitive activities labelled “general components of competence in problem solving”: problem representation, solution strategies, self-monitoring, and explanations. Analytical problem solving in everyday contexts, as measured by the ALL problem-solving instrument, focuses on the components 1 to 3 listed above (and to some extent 4).

One of the most important insights of recent research in cognitive psychology is that solving demanding problems requires at least some knowledge of the domain in question. The concept of a problem space through which a General Problem Solver moves by means of domain-independent search strategies (Newell and Simon, 1972) proved to be too simple to describe how problem situations are understood and the process of finding a solution. Efforts to identify a general, domain-independent competence for steering dynamic systems (operative intelligence) within the framework of complex problem-solving research were also unsuccessful; performance on such systems can only partially be transferred to other systems (Funke, 1991). However, research on grade 3 to grade 12 students showed that problem-solving skills clearly improve under well-tuned training conditions and that a substantial transfer across different problems can be achieved (Reeff et al. 1989, 1992, 1993; Regenwetter, 1992; Regenwetter and Müller, 1992; Stirner, 1993).

Problem solving is dependent on knowledge of concepts and facts (declarative knowledge) and knowledge of rules and strategies (procedural knowledge) in a given subject domain. Although it is evident from past research that declarative knowledge in the problem domain can substantially contribute to successful problem-solving strategies, procedural knowledge is

crucial as well. The amount of relevant previous knowledge available could also account for the relation between intelligence and problem-solving performance, as shown in the work of Raaheim (1988) and Leutner (1999). People with no relevant previous knowledge at all are unable to explore the problem situation or plan a solution in a systematic manner and are forced to rely on trial and error instead. Those who are already very familiar with the task are able to deal with it as a matter of routine. General intellectual ability, as measured by reasoning tasks, plays no role in either of these cases. When problem solvers are moderately familiar with the task, analytical reasoning strategies can be successfully implemented.

The approach taken for the assessment of problem solving in ALL relies on the notion of (moderately) familiar tasks. Within a somewhat familiar context the problems to be solved are inexplicit enough so as not to be perceived as pure routine tasks. On the other hand, the domain-specific knowledge prerequisites are sufficiently limited as to make analytical reasoning techniques the main cognitive tool for solving the problems.

4.5.1 Identifying task characteristics

How can contextualized, real-life problems be defined and transformed into a set of assessment tasks? After reviewing the various approaches that have been taken in previous research to measure problem solving, a decision was made to use a project approach in ALL. The project approach has the potential to be a powerful means for assessing analytical problem solving skills in real world, everyday contexts for several reasons. Solving problems in project-like settings is important and relevant for adults in both their professional and their private life. In addition, the project approach has been successfully implemented in other large-scale assessments, and it can be realized as a paper-and-pencil-instrument, which is of crucial importance for contemporary large-scale surveys. Furthermore, the project approach uses different problem-solving stages as a dimension along which to generate the actual test items. Following Pólya (1945, 1980), the process of problem solving has been frequently described in terms of the following stages:

- ❖ Define the goal.
- ❖ Analyze the given situation and construct a mental representation.
- ❖ Devise a strategy and plan the steps to be taken.
- ❖ Execute the plan, including control and – if necessary – modification of the strategy.
- ❖ Evaluate the result.

The different action steps define the course of action for an “everyday” project. One or more tasks or items are generated to correspond to each of these action steps. Respondents are expected to work on individual tasks that have been identified as steps that need to be carried out as a part of their project (a sample project, for example, might involve “planning a reunion” or “renovating a clubhouse”). Embedding the individual tasks in a project is believed to yield a high degree of context authenticity. Although they are part of a comprehensive and coherent project, the individual tasks are designed so that they can be solved independently of one another and are expected to vary in complexity and overall difficulty for adults.

Since assessing problem solving skills in large-scale assessments is a relatively new endeavour, it might be helpful to provide a detailed account of the construction process. Table A1 provides an overview of the problem solving steps as they correspond to the action steps identified above. Different components and aspects of each of the problem solving steps are listed.

Table A1 Problem-solving steps and instantiations

Define the goals	<ul style="list-style-type: none"> • Set goals. • Recognize which goals are to be reached and specify the essential reasons for the decision. • Recognize which goals/wishes are contradictory and which are compatible. • Assign priorities to goals/wishes.
Analyze the situation	<ul style="list-style-type: none"> • Select, obtain and evaluate information. <ul style="list-style-type: none"> ⇒ What information is required, what is already available, what is still missing, and what is superfluous? ⇒ Where and how can you obtain the information? ⇒ How should you interpret the information? • Identify the people (e.g. with what knowledge and skills) who are to be involved in solving the problem. • Select the tools to be used. • Recognize conditions (e.g. time restrictions) that need to be taken into account.
Plan the solution	<ul style="list-style-type: none"> • Recognize which steps need to be taken. • Decide on the sequence of steps (e.g. items on the agenda). • Coordinate work and deadlines. • Make a comparative analysis of alternative plans (recognize which plan is suitable for reaching the goals). • Adapt the plan to changed conditions. • Opt for a plan.
Execute the plan	<ul style="list-style-type: none"> • Carry out the individual steps (e.g., write a letter, fill in a form, make calculations).
Evaluate the results	<ul style="list-style-type: none"> • Assess whether and to what extent the target has been reached. • Recognize mistakes. • Identify reasons for mistakes. • Assess consequences of mistakes.

The construction of a pool of assessment tasks that could be mapped back to these five action steps involved several phases of activities. First was the identification of appropriate projects that would be suitable for adults with varying educational backgrounds and relevant to the greatest number of people in the target group. Next, developers had to identify and sketch out the problem situation and the sequence of action steps that relate back to the model. Third, they had to develop a pool of items that were consistent with the action steps and that tapped into particular processes including the development of correct responses and appropriate distractors for multiple choice items and solution keys and scoring guides for open-ended tasks.

4.5.2 Characterizing problem solving tasks

ALL included a total of 4 projects involving 20 tasks in the assessment of problem solving. These resulted in 19 scorable items that ranged from 199 to 394 along the scale and, like the literacy and numeracy tasks, their placement was determined by the patterns of right and wrong responses among adults in participating countries. Rather than release one of the four projects

that were used in ALL, we will characterize the hypothesized proficiency scale for analytical problem solving that was tested using pilot data and present an example from the pilot data that was not used in the main assessment³. Similar models have been described within the frameworks of other large-scale assessments of problem-solving competencies such as the project test for Hamburg/Germany (Ebach, Klieme and Hensgen, 2000) and the PISA 2003 assessment of cross-curricular problem solving (OECD, in press).

In ALL, four levels of problem-solving proficiency are postulated:

Level 1

At a very elementary level, concrete, limited tasks can be mastered by applying content-related, practical reasoning. At this level, people will use specific content-related schemata to solve problems.

Level 2

The second level requires at least rudimentary systematical reasoning. Problems at this level are characterized by well-defined, one-dimensional goals; they ask for the evaluation of certain alternatives with regard to transparent, explicitly stated constraints. At this level, people use concrete logical operations.

Level 3

At the third level of problem-solving proficiency, people will be able to use formal operations (e.g., ordering) to integrate multi-dimensional or ill-defined goals, and to cope with non-transparent or multiple dependent constraints.

Level 4

At the final and highest level of competency, people are capable of grasping a system of problem states and possible solutions as a whole. Thus, the consistency of certain criteria, the dependency among multiple sequences of actions and other “meta-features” of a problem situation may be considered systematically. Also, at this stage people are able to explain how and why they arrived at a certain solution. This level of problem-solving competency requires a kind of critical thinking and a certain amount of meta-cognition

The following example illustrates a concrete realization of a project. For this purpose a project that is not included in the final ALL instrument is introduced and one typical problem-solving task is shown. The project is about “Planning a trip and a family reunion”.

In the introductory part of the project, the respondent is given the following summary describing the scenario and overall problem:

“Imagine that you live in City A. Your relatives are scattered throughout the country and you would like to organize a family reunion. The reunion will last 1 day. You decide to meet in City B, which is centrally located and accessible to all. Since you and your relatives love hiking, you decide to plan a long hike in a state park close to City B. You have agreed to be responsible for most of the organization.”

The respondent is then given a list of steps he or she needs to work through, in this example the following list:

- ❖ Set the date for the reunion
- ❖ Consider your relatives' suggestions for the hike
- ❖ Plan what needs to be done before booking your flight
- ❖ Answer your relative's questions about traveling by plane
- ❖ Book your flight
- ❖ Make sure your ticket is correct
- ❖ Plan the trip from City B to the airport

The first task of this project "Set the date for the reunion" is a good example of a typical problem-solving task and is shown here as it would appear in a test booklet.

Example task: Set the date for the reunion

The family reunion should take place sometime in July.

You asked all your relatives to tell you which dates would be suitable. After talking to them, you made a list of your relatives' appointments during the month of July. Your own appointment calendar is lying in front of you. You realize that some of your relatives will have to arrive a day early in order to attend the family reunion and will also only be able to return home on the day after the meeting.

Please look at the list of your relatives' appointments and your own appointment calendar.

List of your relatives' appointments in July 1999

Henry	Karen	Peter	Janet	Anne	Frank
Vacation in City E beginning July 16; Appointment on July 11	Every day of the week is okay except Thursdays and on July 16	Business appointments on July 2, July 13, and between July 27 and 29	Doesn't have any appointments	Unable to attend reunion on July 5, July 20, or July 24	Has to be away sometime during the 1 full week in July on business, but will find out the exact dates shortly before

Henry, Karen, and Peter could arrive on the same day as the reunion whereas Janet, Anne, and Frank can only arrive on the afternoon before and return home on the day after the reunion.

Example task (cont.)

Your appointment calendar for July 1999

July 1999

Thurs.	1	Meeting with David
Fri.	2	
Sat.	3	
Sun.	4	
Mon.	5	
Tue.	6	
Wed.	7	
Thurs.	8	
Fri.	9	
Sat.	10	Hike in City C
Sun.	11	
Mon.	12	
Tue.	13	
Wed.	14	
Thurs.	15	
Fri.	16	
Sat.	17	
Sun.	18	
Mon.	19	
Tue.	20	
Wed.	21	
Thurs.	22	
Fri.	23	
Sat.	24	
Sun.	25	
Mon.	26	
Tue.	27	
Wed.	28	Vacation
Thurs.	29	Vacation
Fri.	30	Vacation
Sat.	31	

Question 1. Which of the following dates are possible for the family reunion?

Please select all possible dates.

a July 4

b	July 7
c	July 14
d	July 18
e	July 25
f	July 29

This project illustrates nicely how the action steps logic is actually “translated” into a concrete thematic action flow. The underlying plot – planning a trip and a family reunion – constitutes a very typical everyday-type of action that presumably a large majority of people in different countries will be able to relate to. The action steps themselves and their sequence can deviate from the normative complete action model, as is the case here. The normative model is used as a guideline that is adapted to each specific context. In this case, for example, the task “Consider your relatives’ suggestions for the hike” corresponds approximately to the action step “Analyze the situation”, the task “Plan what needs to be done before booking your flight” corresponds to the action step “Plan the solution”, and “Book your flight” is a typical example for the action step “Execute the plan”.

The example task gives a first indication of item structures and formats. The tasks typically start off with a short introduction to the situation, followed by varying types and amounts of information that need to be worked through. In the example task, in order to set the date for the family reunion, the respondent needs to process, compare and integrate the information provided in the list of the relatives’ appointments, including the addendum to this list, and their own appointment calendar. Here the information is mostly textual and in the form of tables. The answer format is a multiple-choice format with more than one correct response alternatives, although the number of correct response alternative is not specified.

4.6 Conclusion

This chapter has offered a brief overview of the frameworks that have been used for both developing the tasks used to measure prose and document literacy, numeracy and problem solving in ALL as well as for understanding the meaning of what is being reported with respect to the comparative literacy proficiencies of adults. The frameworks identify a set of variables that have been shown to influence successful performance on a broad array of tasks. Collectively, they provide a means for moving away from interpreting survey results in terms of discrete tasks or a single number, and towards identifying levels of performance sufficiently generalized to have validity across assessments and groups. As concern ceases to center on discrete behaviours or isolated observations and focuses more on providing meaningful interpretations of performance, a higher level of measurement is reached (Messick, 1989).

4.6.1 Some analytical considerations

The skill levels presented in the ALL dataset not only provide a means for exploring the progression of information-processing demands across each of the scales, but also can be used to help explain how the proficiencies individuals demonstrate reflect the likelihood they will respond correctly to the broad range of tasks used in this assessment as well as to any task that has the same characteristics. In practical terms, this means that individuals performing at 250 on each scale are expected to be able to perform the average level 1 and 2 tasks with a high degree of proficiency – i.e. with an average probability of a correct response at 80 per cent or higher. It does not mean that they will not be able to perform tasks in levels 3 or higher. They would be expected to do so some of the time, but not consistently.

Based on the result of the 1994 IALS for the two scales common to the 2003 ALL, Tables 4.1 and 4.2 display the probability that individuals performing at selected points on the prose or document literacy scales will give a correct response to tasks of varying difficulty. For example, a reader whose prose proficiency is 150 has less than a 50 per cent chance of giving a correct response to the level 1 tasks. Individuals whose proficiency score is 200, in contrast, have about an 80 per cent probability of responding correctly to these tasks.

In terms of task demands, it can be inferred that adults performing at 200 on the prose scale are likely to be able to locate a single piece of information in a brief text when there is no distracting information, or if plausible but incorrect information is present but located away from the correct answer. However, these individuals are likely to encounter far more difficulty with tasks in levels 2 through 5. For example, they would have only a 40 per cent chance of performing the average level 2 task correctly, an 18 per cent chance of success with tasks in level 3, and no more than a 7 per cent chance with tasks in levels 4 and 5.

In contrast, respondents demonstrating a proficiency of 300 on the prose scale have about an 80 per cent chance or higher of succeeding with tasks in levels 1, 2 and 3. This means that they demonstrate success with tasks that require them to make low-level inferences and with those that entail taking some conditional information into account. They can also integrate or compare and contrast information that is easily identified in the text. On the other hand, they are likely to encounter difficulty with tasks where they must make more sophisticated text-based inferences, or where they need to process more abstract types of information. These more difficult tasks may also require them to draw on less familiar or more specialised types of knowledge beyond that given in the text. On average, they have about a 50 per cent probability of performing level 4 tasks correctly; with level 5 tasks, their likelihood of responding correctly decreases to 40 per cent.

Similar kinds of interpretations can be made using the information presented for the document and quantitative literacy scales. For example, someone who is at 200 on the quantitative scale has, on average, a 67 per cent chance of responding correctly to level 1 tasks. His or her likelihood of responding correctly decreases to 47 per cent for level 2 tasks, 21 per cent for

Where have all the illiterates gone?

Like its predecessor, the 2003 ALL conceptualizes proficiency along a continuum that denotes how well adults use information to function in society and the economy. It bares repeating that the ALL does not measure the absence of competencies. Rather it measures knowledge and skills in the four domains along a broad range of ability. Consequently, the results cannot be used to classify population groups as either “literate” or “illiterate”.

level 3 tasks, 6 per cent for level 4 tasks and a mere 2 per cent for level 5 tasks. Similarly, readers with a proficiency of 300 on the quantitative scale would have a probability of 92 per cent or higher of responding correctly to tasks in levels 1 and 2. Their average probability would decrease to 81 per cent for level 3 tasks, 57 per cent for level 4 and 20 per cent for level 5.

Table 4.1 **Average probabilities of successful performance, prose scale**

Prose level	Selected proficiency scores				
	150	200	250	300	350
			%		
1	48	81	95	99	100
2	14	40	76	94	99
3	6	18	46	78	93
4	2	7	21	50	80
5*	2	6	18	40	68

* Based on one task

Source: Adult Literacy Survey (1994).

Table 4.2 **Average probabilities of successful performance, document scale**

Document Level	Selected proficiency scores				
	150	200	250	300	350
			%		
1	40	72	94	99	100
2	20	51	82	95	99
3	7	21	50	80	94
4	4	13	34	64	85
5*	<1	1	3	13	41

* Based on one task

Source: Adult Literacy Survey (1994).

Proficiency in each domain is measured on a continuous scale. Each scale starts at zero and increases to a theoretical maximum of 500 points (with four decimal places of precision). Scores along the scale denote the points at which a person with a given level of performance has an 80 percent probability of successfully completing a task at that level of difficulty.

From an analytical standpoint, these continuous measures can be useful in creating summary statistics that describe the competencies of populations such as their overall average score. Populations with similar average scores, however, may have quite different numbers of low or high performing adults. Thus, one can also look at how the scores are distributed within populations by using percentile scores. Percentile scores are the scores below which a specified percentage of adults are found. Thus, for example, the 5th percentile score is the one below which we find 5% of adults in a particular population. Differences in percentile scores tell us something about the degree of equality in proficiency across populations. Users should refer to chapter 5 for more detailed information on the correct use of the plausible values in their analysis.

The ALL scores have also been grouped into proficiency levels representing tasks of increasing difficulty. For the prose and document literacy domains as well as the numeracy domain, experts have defined five broad levels of difficulty, each corresponding to a similar (though not equidistant) range of scores. For the problem solving domain, experts have defined four broad levels of difficulty. In each domain, Level 1 denotes the group with the lowest proficiency and Level 4 for the problem solving and 5 for the other domains contains the highest performers.

It is important, for analytical as well as operational reasons, to define a “desired level” of competence for coping with the increasing skill demands of the emerging knowledge and information economy. Level 3 performance is generally chosen as a benchmark because in developed countries, performance above Level 2 is generally associated with a number of positive outcomes. These include increased civic participation, increased economic success and independence, and enhanced opportunities for lifelong learning and personal literacy (Kirsch, I., et al., 1993; Murray, T.S. et al., 1997; Tuijnman, A., 2001). Whereas individuals at proficiency Levels 1 and 2 typically have not yet mastered the minimum foundation of literacy needed to attain higher levels of performance (Strucker, J., Yamamoto, K. 2005)


Secondary analysis of the 1994 IALS data has yielded consistent evidence that the performance difference between Level 2 and Level 3 on the prose, document and quantitative literacy scales is substantive and corresponds to a significant difference in measurable benefits accruing to citizens in OECD countries (OECD and HRDC, 1997). Results of preliminary analysis of the ALL data, including the new numeracy scale, are consistent with this finding. For this reason, it is sometimes useful to anchor the scales at the cut point between Levels 2 and 3, thus highlighting the distributions above and below this threshold for the prose, document and numeracy domains. In contrast, interpretation of the problem solving domain is more complex and no single “desirable” threshold has yet been set, in which case, a cutpoint at level 1 may be more appropriate until a more precise threshold can be found.

Chapter 8 will offer further tools and techniques for the appropriate use of the Plausible Values and Replicate Weights required to produce accurate estimates of the standards errors associated with each point estimate.

References

- Almond, R.G., and Mislevy, R.J. (1998). Graphical models and computerized adaptive testing. (TOEFL Tech. Rep. No. 14). Princeton, NJ: Educational Testing Service.
- Baker, D., and Street, B. (1994). Literacy and numeracy: Concepts and definitions. In T. Husen and E. A. Postlethwaite (Eds.), *Encyclopedia of education*. New York: Pergamon Press.
- Beazley, K. (1984). *Education in Western Australia: Report of the Committee of Inquiry into Education in Western Australia*. Education Department of Western Australia.
- Coben, D., O'Donoghue, J., and FitzSimons, G.E. (Eds.)(2000). *Perspectives on adults learning mathematics: Theory and practice*. London: Kluwer Academic Publishers.
- Cockcroft, W.H. (1982). *Report of the Committee of Inquiry into the Teaching of Mathematics in Schools*. London: HMSO.
- Cook-Gumperz, J., and Gumperz, J. (1981). From oral to written culture: The transition to literacy. In M. Whitman (Ed.), *Writing: The nature, development and teaching of written communication: Vol. 1*. Hillsdale, NJ: Erlbaum.
- Crandall, J. (1981, December). Functional literacy of clerical workers: Strategies for minimizing literacy demands and maximizing available information. Paper presented at the annual meeting of the American Association for Applied Linguistics, New York.
- Diehl, W. (1980). *Functional literacy as a variable construct: An examination of the attitudes, behaviours, and strategies related to occupational literacy*. Unpublished doctoral dissertation, Indiana University.
- Dossey, J.A. (1997). "Defining and measuring quantitative literacy". In L.A. Steen (Ed.), *Why numbers count: Quantitative literacy for tomorrow's America*. New York: College Entrance Examination Board.
- Fey, James T. (1990). "Quantity" In L.A. Steen (Ed.) *On the shoulders of giants: New approaches to numeracy*. Washington, DC: National Academy Press.
- Frankenstein, M. (1989). *Relearning mathematics: A different third 'R' – Radical maths*. London: Free Association Books.
- Gal, I. (1997). Numeracy: Imperatives of a forgotten goal. In L.A. Steen (Ed.), *Why numbers count: quantitative literacy for tomorrow's America* (pp. 36-44). New York: The College Board.
- Gal, I. (2000). The numeracy challenge. In I. Gal (Ed.), *Adult numeracy development: Theory, research, practice* (pp. 1-25). Cresskill, NJ: Hampton Press.
- Gal, I. (2002). Adult Statistical literacy: Meanings, components, responsibilities. *International Statistical Review*, 70(1), 1-25.
- Jacob, E. (1982). *Literacy on the job: Final report of the ethnographic component of the industrial literacy project*. Washington, DC: Center for Applied Linguistics.
- Johnston, B. (1994, Summer). Critical numeracy? In *Fine print*, Vol. 16, No. 4.
- Heath, S.B. (1980). The functions and uses of literacy. *Journal of Communication*, 30, 123–133.
- Kirsch, I.S., and Guthrie, J.T. (1984a). Adult reading practices for work and leisure. *Adult Education Quarterly*, 34(4), 213–232.
- Kirsch, I.S., and Guthrie, J.T. (1984b). Prose comprehension and text search as a function of reading volume. *Reading Research Quarterly*, 19, 331–342.

- Kirsch, I. (2001). The International Adult Literacy Survey (IALS): Understanding What Was Measured (ETS Research Report RR-01-25). Princeton, NJ: Educational Testing Service.
- Marr, B., and Tout, D. (1997). A numeracy curriculum: Australian Association of Mathematics Teachers (AAMT) conference proceedings. Melbourne: AAMT.
- Messick, S. (1989). Validity. In R. Linn (Ed.), Educational measurement (3rd ed.). New York: Macmillan.
- Messick, S. (1994). The interplay of evidence and consequences in the validation of performance assessments. *Education Researcher*, 32(2), 13-23.
- Mikulecky, L. (1982). Job literacy: The relationship between school preparation and workplace actuality. *Reading Research Quarterly*, 17(3), 400-419.
- Miller, P. (1982). Reading demands in a high-technology industry. *Journal of Reading*, 26(2), 109-115.
- Mislevy, R.J. (September, 2000). Leverage points for improving educational assessment. Paper submitted to National Center for Research on Evaluation, Standards, and Student Testing (CRESST) as part of award #R305B60002 from the US Department of Education, Office of Educational Research and Improvement.
- Mosenthal, P.B., and Kirsch, I.S. (1998). A new measure for assessing document complexity: The PMOSE/IKIRSCH document readability formula. *Journal of Adolescent and Adult Literacy*, 41(8), 638-657.
- Murray, T.S., Clermont, Y. and Binkley, M. (Eds.) The Adult Literacy and Life Skills Survey: Aspects of Design, Development and Validation. Canada: Statistics Canada, in press.
- Murray, T.S., Kirsch, I.S., and Jenkins, L. (1998). Adult Literacy in OECD Countries: Technical report on the First International Adult Literacy Survey. Washington, DC: National Center for Education Statistics.
- National Council of Teachers of Mathematics. (2000). Principles and standards for school mathematics. Reston, VA: Author.
- Organization for Economic Co-operation and Development. (1992). Adult illiteracy and economic performance. Paris, France: Author.
- Rutherford, F.J and Ahlgren, A. (1990). Science for all Americans. New York: Oxford University Press.
- Rychen, D.S. and Salganik, L.H. (Eds.) Key Competencies for a Successful Life and a Well-Functioning Society. Cambridge, MA: Hogrefe and Huber Publishers, 2003.
- Senechal, Majorie (1990) "Shape" In L.A. Steen (Ed.) On the shoulders of giants: New approaches to numeracy. Washington, DC: National Academy Press.
- Scribner, S., and Cole, M. (1981). The psychology of literacy. Cambridge, MA: Harvard University Press.
- Steen, L.A. (Ed). (1990). On the shoulders of giants: New approaches to numeracy. Washington, DC: National Research Council.
- Steen, L.A. (2001). Mathematics and democracy: the case for quantitative literacy. USA: National Council on Education and the Disciplines.
- Sticht, T.G. (Ed.). (1975). Reading for working: A functional literacy anthology. Alexandria, VA: Human Resources Research Organization.
- Sticht, T.G. (1978). Literacy and vocational competency (Occasional Paper 39, National Center for Research in Vocational Education). Columbus, OH: Ohio State University.

- 
- Sticht, T.G. (1982, January). Evaluation of the reading potential concept for marginally literate adults. (Final Report FR-ET50-82-2). Alexandria, VA: Human Resources Research Organization.
- Szwed, J. (1981). The ethnography of literacy. In M. Whitman (Ed.), *Writing: The nature, development, and teaching of written communication*: Vol. 1. Hillsdale, NJ: Erlbaum.
- Tobias, S. (1993). *Overcoming math anxiety*. New York: Norton.
- Venezky, R.L. (1983). The origins of the present-day chasm between adult literacy needs and school literacy instruction. *Visible Language*, 16, 113-136.

5.0 Survey Methodology

Each participating country was required to design and implement the Adult Literacy and Life Skills (ALL) survey according to the standards provided in the document '*Standards and Guidelines for the Design and Implementation of the Adult Literacy and Life Skills Survey*'. These ALL standards established the minimum survey design and implementation requirements for the following project areas:

1. Survey planning	12. Respondent contact strategy
2. Target population	13. Response rate strategy
3. Method of data collection	14. Interviewer hiring, training, supervision
4.. Sample frame	15. Data capture
5. Sample design	16. Coding
6. Sample selection	17. Scoring
7. Literacy assessment design	18. All data file-format and editing
8. Background questionnaire	19. Weighting
9. Task booklets	20. Estimation
10. Instrument requirements to facilitate data processing	21. Confidentiality
11. Data collection	22. Survey documentation
	23. Pilot Survey

5.1 Assessment design

The participating countries, with the exception of the state of Nuevo Leon in Mexico, implemented an ALL assessment design. Nuevo Leon assessed literacy using the International Adult Literacy Survey (IALS) assessment instruments.

In both ALL and IALS a Balanced Incomplete Block (BIB) assessment design was used to measure the skill domains. The BIB design comprised a set of assessment tasks organized into smaller sets of tasks, or blocks. Each block contained assessment items from one of the skill domains and covers a wide range of difficulty, i.e., from easy to difficult. The blocks of items were organized into task booklets according to a BIB design. Individual respondents were not required to take the entire set of tasks. Instead, each respondent was randomly administered one of the task booklets.

ALL assessment

The ALL psychometric assessment consisted of the domains Prose, Document, Numeracy, and Problem Solving. The assessment included four 30-minute blocks of Literacy items (i.e., Prose AND Document Literacy), two 30-minute blocks of Numeracy items, and two 30-minute blocks of Problem-Solving items.

A four-domain ALL assessment was implemented in Bermuda, Canada, Hungary, Italy, Netherlands, New Zealand, Norway, and the French and German language regions of

Switzerland. The United States and the Switzerland Italian language region carried out a three-domain ALL assessment that excluded the Problem Solving domain. In addition to the mentioned assessment domains, these participating countries assessed the use of information and communication technology via survey questions incorporated in the ALL Background Questionnaire.

The blocks of assessment items were organized into 28 task booklets in the case of the four-domain assessment and into 18 task booklets for the three domain assessment. The assessment blocks were distributed to the task booklets according to a BIB design whereby each task booklet contained two blocks of items. The task booklets were randomly distributed amongst the selected sample. In addition, the data collection activity was closely monitored in order to obtain approximately the same number of complete cases for each task booklet, except for two task booklets in the three-domain assessment containing only Numeracy items that required a larger number of complete cases.

IALS assessment

The state of Nuevo Leon, Mexico carried out an IALS assessment. The IALS assessment consisted of three literacy domains: Prose, Document, and Quantitative. In addition, the ALL Background Questionnaire was used in Nuevo Leon. The use of information and communication technology was assessed via survey questions incorporated in the ALL Background Questionnaire.

IALS employed seven task booklets with three blocks of items per booklet. The task booklets were randomly distributed amongst the selected sample. In addition, the data collection activity was monitored in order to obtain approximately the same number of complete cases for each task booklet.

5.2 Target population and sample frame

Each participating country designed a sample to be representative of its *civilian non-institutionalized persons 16 to 65 years old (inclusive)*.

Countries were also at liberty to include adults over the age of 65 in the sample provided that a minimum suggested sample size requirement was satisfied for the 16 to 65 year age group. Canada opted to include in its target population adults over the age of 65. All remaining countries restricted the target population to the 16 to 65 age group.

Exclusions from the target population for practical operational reasons were acceptable provided a country's survey population did not differ from the target population by more than five percent, i.e. provided the total number of exclusions from the target population due to undercoverage was not more than five percent of the target population. All countries indicate that this five-percent requirement was satisfied.

Each country chose or developed a sample frame to cover the target population. The following table shows the sample frame and the target population exclusions for each country:

TABLE 5.1 Sample frame and target population exclusions

Country	Sample frame	Exclusions
Bermuda	Land Valuation List □ an up-to-date listing of all housing units in Bermuda.	Persons residing in institutions, visitors to Bermuda (i.e., persons staying less than 6 months).
Canada	Census of Population and Housing database, reference date of May 15, 2001 □ households enumerated by the Census long-form (20% sample)	Long-term institutional residents, members of the armed forces, individuals living on Indian Reserves, residents of sparsely populated regions.
Hungary	Census of Population and Housing database.	Homeless people, prisoners.
Italy	Polling list – a list of individuals aged 18 and over that are resident in Italy and have civil rights	None
Netherlands	Municipal Basic Administration (GBA) as collected by the National Statistical Office (CBS) on a monthly basis.	Persons living in institutions; persons illegally in the country.
New Zealand	Census Meshblocks (as developed for the New Zealand Census by Statistics New Zealand)	Persons living in non-private dwellings such as prisons, retirement homes, hospitals, university residences etc; persons living in remote rural areas and on off-shore islands (except Waiheke Island which is included).
Norway	Norwegian Register of Education (2002 version)	Permanent residents in institutions, individuals for whom education level is unknown
Nuevo Leon, Mexico	Census of Population and Housing database, reference year 2000	Persons residing in institutions, members of the Mexican Navy
Switzerland	Register of private telephone numbers (September 2002)	Persons living in institutions, people living in very isolated areas, persons with no private telephone number
United States	Area Frame – 1,883 Primary Sampling Units covering all counties in the 50 states in the United States plus Washington, DC	Full-time military personnel, residents in institutionalized group quarters

5.3 Sample design

Each participating country was required to use a probability sample representative of the national population aged 16 to 65. Of course, the available sampling frames and resources varied from one country to another. Therefore, the particular probability sample design to be used was left to the discretion of each country. Each country's proposed sample design was reviewed by Statistics Canada to ensure that the sample design standards and guidelines were satisfied.

Each country's sample design is summarized below. The sample size and response rate for each country can be found in the section following this one.

Bermuda

A two-stage stratified probability design was employed. In stage one Bermuda's Land Valuation List of dwellings was stratified by parish, i.e., geographic region. Within each parish, a random sample of dwellings was selected with probability proportional to the number of parish dwellings. At stage two, one eligible respondent was selected using a Kish-type person selection grid.

Canada

A stratified multi-stage probability sample design was used to select the sample from the Census Frame. The sample was designed to yield separate samples for the two Canadian official languages, English and French. In addition, Canada increased the sample size in order to produce estimates for a number of population subgroups. Provincial ministries and other organizations sponsored supplementary samples to increase the base or to target specific subpopulations such as youth (ages 16 to 24 in Québec and 16 to 29 in British Columbia), adults aged 25 to 64 in Québec, linguistic minorities (English in Québec and French elsewhere), recent and established immigrants, urban aboriginals, and residents of the northern territories.

In each of Canada's ten provinces the Census Frame was further stratified into an urban stratum and a rural stratum. The urban stratum was restricted to urban centers of a particular size, as determined from the previous census. The remainder of the survey frame was delineated into primary sampling units (PSUs) by Statistics Canada's Generalised Area Delineation System (GARDS). The PSUs were created to contain a sufficient population in terms of the number of dwellings within a limited area of reasonable compactness. In addition, the Census Frame was ordered within each geographic region by highest level of education prior to sample selection, thus ensuring a representation across the range of educational backgrounds.

Within the urban stratum, two stages of sampling were used. In the first stage, households were selected systematically with probability proportional to size. During the second stage, a simple random sample algorithm was used by the CAPI application to select an individual from the eligible household adults. Three stages were used to select the sample in the rural stratum. In the first stage, Primary Sampling Units were selected with probability proportional to population size. The second and third stages for the rural stratum repeated the same methodology employed in the two-stage selection for the urban stratum.

Hungary

A stratified two-stage sample design was employed to yield a sample of persons selected Proportional to Population Size (PPS).

The population was stratified into seven regions and twenty counties. This stratification took into consideration the regional and county demographic characteristics and other conditions (e.g. rate of active and inactive population, unemployment rate) that varied from one region to another. In each county, the population was further stratified into three types of settlements: city, town, and village. Subsequently, the sample was selected in two stages:

Stage1: a PPS sample of settlements,

Stage2: a random selection of addresses from the settlements selected at stage 1. The list of addresses in each selected settlement were obtained from the Ministry of Interior files from the 2001 Census, the most up-to-date and precise data for the population of Hungary at the time of sample selection. The addresses to be contacted for interview were selected from these files.

Italy

The sample was selected from the 2002 version of the Norwegian Register of Education using a two-stage probability sample design.

The design created 363 primary sampling units (PSUs) from the 435 municipalities in Norway. These PSUs were grouped into 109 geographical strata. Thirty-eight strata consisted of one PSU that was a municipality with a population of 25,000 or more. At the first stage of sample selection, each of these 38 PSUs was included with certainty in the sample. The remaining municipalities were allocated to 79 strata. The variables used for stratification of these municipalities were industrial structure, number of inhabitants, centrality, communication structures, commuting patterns, trade areas and (local) media coverage. One PSU was selected with probability proportional to size from each of these 79 strata.

The second stage of the sample design involved the selection of a sample of individuals from each sampled PSU. Each selected PSU was stratified by three education levels defined by the Education Register. The sample size for each selected PSU was determined by allocating the overall sample size to each selected PSU with probability proportional to the target population size. The PSU sample was then allocated with 30 percent from the low-education group, 40 percent from the medium-education group and 30 percent from the high-education group. Individuals for whom the education level (84,318 persons) was not on the Education Register were excluded from the sampling.

Netherlands

The sample design in the Netherlands was a stratified, multi-stage systematic cluster design.

In the first stage, the country was stratified into 4 regions; North, East, West, and South. Within each stratum, a sample of municipalities was selected with probability proportional to municipality population size. This was achieved by ordering the municipalities within a stratum by population size and by systematically selecting the sample of municipalities using a random starting point and a fixed sampling interval. The population data were based on the municipality data, Gemeentelijke Basis Administratie (GBA), from the national statistical office, Centraal Bureau voor Statistiek (CBS).

In the second stage, within each selected municipality a systematic sample of postal code areas was drawn. The company, Experian, provided information about credit score (i.e. the percentage of households having debts within a postal code area) and purchasing power for the

postal code areas (6 digits). The postal code areas were ordered by credit score and then by purchasing power. From a random starting point and with a fixed sampling interval (in terms of households), the households were drawn.

In the third stage within each selected postal code area one household was randomly selected. Data came from the Experian database on household information (based on CENDRIS, the current owner of the Post Office central database). This database is updated on a monthly basis.

In the fourth stage one eligible individual within the selected household was randomly selected.

New Zealand

The sample design was a stratified probability design with three stages of sampling – replicate, dwelling, and household member. The population was categorized into three strata - main stratum (everyone 16 to 65 eligible), Maori and Pacific stratum (only Maori and Pacific eligible), and Pacific stratum (only Pacific people eligible).

(a) Stage 1: The Replicate

From the 38,000 meshblocks which formed the basis of New Zealand's 2001 Census of Population and Dwellings, those with 9 or fewer dwellings were eliminated, leaving 32,115 meshblocks with 10 or more dwellings. The coverage of permanent private dwellings was 98.6 per cent. The probability of selection for each meshblock was proportional to the number of dwellings in the meshblock. A total of 896 meshblocks were selected, and subsequently allocated to 32 replicates made up of 28 meshblocks per replicate. Each replicate contained meshblocks distributed north to south in approximately the same manner, and was thus a mini national probability sample.

(b) Stage 2: The Dwelling

For the main stratum, dwellings were selected as follows. The sample interval was derived for each meshblock as the number of dwellings in the meshblock divided by 15. The sample interval thus differed according to the size of the meshblock. Beginning from a randomised starting point, interviewers selected dwellings according to the meshblock's sample interval.

In addition to the dwellings in the main stratum, up to an additional 21 dwellings per meshblock were also sampled for the Māori and Pacific, and the Pacific strata. In 4 of these dwellings, residents of either Māori or Pacific ethnicity were eligible for selection. In the remaining 17 dwellings only residents of Pacific ethnicity were eligible. The sample interval was 1 for these dwellings once the main stratum dwellings were set aside.

(c) Stage 3: The Respondent

For the main stratum, one person per household was selected from all eligible household members using a Kish grid. For the two ethnic strata, the ethnicity of the household members (Māori or Pacific for stratum two, Pacific for stratum three) was an additional eligibility criterion prior to selection using the Kish grid.

Nuevo Leon, Mexico

The sample design was a stratified probability design with two stages of sampling within each stratum.

The 51 municipalities in Nuevo Leon were grouped geographically into three strata: Stratum 1 – Census Metropolitan Area of Monterrey, consisting of 9 municipalities; Stratum 2 – the municipalities of Linares and Sabinas Hidalgo; Stratum 3 – the remaining 40 municipalities of Nuevo Leon. The initial sample was allocated to the three strata proportional to the number of dwellings in each stratum.

At the first stage of sample selection, in each stratum a simple random sample of households was selected. The second sampling stage consisted of selecting one person belonging to the target population from each selected household using a Kish-type person selection grid.

Switzerland

The sample design was a stratified probability design with two stages of sampling. Separate estimates were required for Switzerland's three language regions (i.e., German, French, Italian). Thus, the three language regions are the primary strata. Within the language regions, the population was further stratified into the metropolitan areas represented by the cantons of Geneva and Zurich and the rest of the language regions. At the first stage of sampling, in each stratum a systematic sample of households was drawn from a list of private telephone numbers. In the second stage, a single person belonging to the target population was selected from each household using a Kish-type person selection grid.

United States

A stratified multi-stage probability sample design was employed in the United States.

The first stage of sampling consisted of selecting a sample of 60 primary sampling units (PSUs) from a total of 1,883 PSUs that were formed using a single county or a group of contiguous counties, depending on the population size and the area covered by a county or counties. The PSUs were stratified on the basis of the social and economic characteristics of the population, as reported in the 2000 Census. The following characteristics were used to stratify the PSUs: region of the country, whether or not the PSU is a Metropolitan Statistical Area (MSA), population size, percentage of African-American residents, percentage of Hispanic residents, and per capita income. The largest PSUs in terms of a population size cut-off were included in the sample with certainty. For the remaining PSUs, one PSU per stratum was selected with probability proportional to the population size.

At the second sampling stage, a total of 505 geographic segments were systematically selected with probability proportionate to population size from the sampled PSUs. Segments consist of area blocks (as defined by Census 2000) or combinations of two or more nearby blocks. They were formed to satisfy criteria based on population size and geographic proximity.

The third stage of sampling involved the listing of the dwellings in the selected segments, and the subsequent selection of a random sample of dwellings. An equal number of dwellings was selected from each sampled segment.

At the fourth and final stage of sampling, one eligible person was randomly selected within households with fewer than four eligible adults. In households with four or more eligible persons, two adults were randomly selected.

5.4 Sample size

A sample size of 5,400 completed cases in each official language was recommended for each country that was implementing the full ALL psychometric assessment (i.e., comprising the domains Prose and Document Literacy, Numeracy, and Problem-Solving). A sample size of 3,420 complete cases in each official language was recommended if the Problem Solving domain was excluded from the ALL assessment.

A sample size of 3,000 complete cases was recommended for the state of Nuevo Leon, Mexico, which assessed literacy skills with the psychometric task booklets of the International Adult Literacy Survey (IALS).

Table 5.2 shows the final number of respondents (complete cases) for each participating country's assessment language(s). **TABLE 5.2 Sample size by assessment language**

Country	Assessment language	Assessment domains ¹	Number of respondents ²
Bermuda	English	P, D, N, PS	2,696
Canada	English	P, D, N, PS	15,694
	French	P, D, N, PS	4,365
Hungary	Hungarian	P, D, N, PS	5,635
Italy	Italian	P, D, N, PS	6,853
Netherlands	Dutch	P, D, N, PS	5,617
New Zealand	English	P, D, N, PS	7,131
Norway	Bokmal	P, D, N, PS	5,411
Nuevo Leon, Mexico	Spanish	P, D, Q	4,786
Switzerland	French	P, D, N, PS	1,765
	German	P, D, N, PS	1,892
	Italian	P, D, N	1,463
United States	English	P, D, N	3,420

1. P – Prose, D – Document, N – Numeracy, PS – Problem Solving, Q – Quantitative.

2. A respondent's data is considered complete for the purposes of the scaling of a country's psychometric assessment data provided that at least the Background Questionnaire variables for age, gender and education have been completed.

5.5 Data collection

The ALL survey design combined educational testing techniques with those of household survey research to measure literacy and provide the information necessary to make these measures meaningful. The respondents were first asked a series of questions to obtain

background and demographic information on educational attainment, literacy practices at home and at work, labour force information, information communications technology uses, adult education participation and literacy self-assessment.

Once the background questionnaire had been completed, the interviewer presented a booklet containing six simple tasks (Core task). Respondents who passed the Core tasks were given a much larger variety of tasks, drawn from a pool of items grouped into blocks, each booklet contained 2 blocks which represented about 45 items. No time limit was imposed on respondents, and they were urged to try each item in their booklet. Respondents were given a maximum leeway to demonstrate their skill levels, even if their measured skills were minimal.

Data collection for the ALL project took place between the fall of 2003 and early spring of 2008, depending on the country. Table 5.3 presents the collection periods for each participating country.

TABLE 5.3 Survey collection period

Country	Collection date
Bermuda	March through August 2003
Canada	March through September 2003
Hungary	July 2007 through February 2008
Italy	May 2003 through January 2004
Netherlands	July 2007 through January 2008
New Zealand	August 2005 - April 2007
Norway	January through November 2003
Nuevo Leon, Mexico	October 2002 through March 2003
Switzerland	January through November 2003
United States	January through June 2003

To ensure high quality data, the ALL Survey Administration Guidelines specified that each country should work with a reputable data collection agency or firm, preferably one with its own professional, experienced interviewers. The manner in which these interviewers were paid should encourage maximum response. The interviews were conducted in home in a neutral, non-pressured manner. Interviewer training and supervision was to be provided, emphasizing the selection of one person per household (if applicable), the selection of one of the 28 main task booklets (if applicable), the scoring of the core task booklet, and the assignment of status codes. Finally the interviewers' work was to have been supervised by using frequent quality checks at the beginning of data collection, fewer quality checks throughout collection and having help available to interviewers during the data collection period.

The ALL took several precautions against non-response bias, as specified in the ALL Administration Guidelines. Interviewers were specifically instructed to return several times to non-respondent households in order to obtain as many responses as possible. In addition, all countries were asked to ensure address information provided to interviewers was as complete as possible, in order to reduce potential household identification problems.

Countries were asked to complete a debriefing questionnaire after the Main study in order to demonstrate that the guidelines had been followed, as well as to identify any collection problems they had encountered. Table 5.4 presents information about interviews derived from this questionnaire.

TABLE 5.4 Interviewer information

Country	Number of languages	Number of interviewers	Average assignment size	Interviewer experience
Bermuda	1	105	40	No specific information provided.
Canada	2	317	62	Professional interviewers with at least 2 years experience.
Hungary	1	175	32	Professional interviewers with at least 2 years experience.
Italy	1	150	45	Professional interviewers, most of which had at least 2 years experience.
Netherlands	1	277	35	Professional interviewers, approximately one fifth of them had no previous survey experience.
New Zealand	1	160	45	Professional interviewers, but interviewer experience not recorded.
Norway	1	320	30	Only a third of the interviewers had at least 2 years experience, the others were trained specifically for this survey.
Nuevo Leon, Mexico	1	209	29	Approximately 70% of interviewers had 2 years of experience.
Switzerland	3	110	60	No specific information provided.
United States	1	106	64	Professional interviewers approximately a quarter of which had no previous survey experience.

5.6 Data Processing

As a condition of their participation in the ALL study, countries were required to capture and process their files using procedures that ensured logical consistency and acceptable levels of data capture error. Specifically, countries were advised to conduct complete verification of the

captured scores (i.e. enter each record twice) in order to minimize error rates. Because the process of accurately capturing the task scores is essential to high data quality, 100 per cent keystroke verification was required.

Each country was also responsible for coding industry, occupation, and education using standard coding schemes such as the International Standard Industrial Classification (ISIC), the International Standard Classification for Occupation (ISCO) and the International Standard Classification for Education (ISCED). Coding schemes were provided by Statistics Canada for all open-ended items, and countries were given specifics instructions about coding of such items.

In order to facilitate comparability in data analysis, each ALL country was required to map its national dataset into a highly structured, standardized record layout. In addition to specifying the position, format and length of each field, the international record layout included a description of each variable and indicated the categories and codes to be provided for that variable. Upon receiving a country's file, Statistics Canada performed a series of range checks to ensure compliance to the prescribed format, flow and consistency edits were also run on the file. When anomalies were detected, countries were notified of the problem and were asked to submit cleaned files.

5.7 Scoring of tasks

Persons charged with scoring in each country received intense training in scoring responses to the open-ended items using the ALL scoring manual. As well they were provided a tool for capturing closed format questions. To aid in maintaining scoring accuracy and comparability between countries, the ALL survey introduced the use of an electronic bulletin board, where countries could post their scoring questions and receive scoring decisions from the domain experts. This information could be seen by all countries who could then adjust their scoring.

To further ensure quality, countries were monitored as to the quality of their scoring in two ways.

First, within a country, at least 20 per cent of the tasks had to be re-scored. Guidelines for intra-country rescoring involved rescoring a larger portion of booklets at the beginning of the scoring process to identify and rectify as many scoring problems as possible. As a second phase, they were to select a smaller portion of the next third of the scoring booklets; the last phase was viewed as a quality monitoring measure, which involved rescoring a smaller portion of booklets regularly to the end of the re-scoring activities. The two sets of scores needed to match with at least 95 percent accuracy before the next step of processing could begin. In fact, most of the intra-country scoring reliabilities were above 95 per cent. Where errors occurred, a country was required to go back to the booklets and rescore all the questions with problems and all the tasks that belonged to a problem scorer.

Second, an international re-score was performed. Each country had 10 per cent of its sample re-scored by scorers in another country. For example, a sample of task booklets from the United States was re-scored by the persons who had scored Canadian English booklets, and vice-versa. The main goal of the re-score was to verify that no country scored consistently differently from another. Inter-country score reliabilities were calculated by Statistics Canada and the results were evaluated by the Educational Testing Service based in Princeton. Again, strict accuracy was demanded: a 90 per cent correspondence was required before the scores were deemed acceptable. Any problems detected had to be re-scored. Table 5.5 shows the high level of inter-country score agreement that was achieved.

TABLE 5.5 Scoring – per cent reliability by domain

Psychometric domain				
Country pairing (rescoring country – original country)	Prose and document (%)	Numeracy (%)	Problem solving (%)	Total (%)
Canada English – Canada French	95	95	92	95
Canada French – Canada English	95	97	94	95
Norway – Canada	91	93	91	92
Canada – United States	94	97	...	95
United States – Canada	95	97	...	95
United States – Bermuda	91	94	...	90
Bermuda – United States	93	95	...	93
Canada French – Switzerland	95	98	97	96
Switzerland – Canada French	94	96	94	95
Switzerland – Italy	96	98	96	96
Italy – Switzerland	93	97	93	94
Canada – Bermuda	83	83
Canada – Nuevo Leon, Mexico	91	95 ¹	...	92
Hungary	94	96	93	94
Netherlands	91	93	93	92
New Zealand	96	97	94	96

... Not applicable.

1. Quantitative literacy.

TABLE 5.6 Scoring operations summary

Country	Scoring start ¹	Number of scorers	Average scoring time per booklet
Bermuda	middle	5	20 min.
Canada	middle	18 ²	13 min.
Hungary	middle	9	20 min.
Italy	beginning	9	15 min.
Netherlands	middle	7	12 min.
New Zealand	beginning	12	20 min.
Norway	middle	17	8 min.
Nuevo Leon, Mexico	middle	12	N.A.
Switzerland	beginning	11	22 min.
United States	beginning	7	12 min.

1. Indicates that the scoring started at the beginning, middle or end of collection.

2. Includes 15 scorers, 2 people to capture problem solving closed format questions and 1 person to capture scoring sheets.

5.8 Survey response and weighting

Each participating country in ALL used a multi-stage probability sample design with stratification and unequal probabilities of respondent selection. Furthermore, there is a need to compensate for the non-response that occurred at varying levels. Therefore, the estimation of population parameters and the associated standard errors is dependent on the survey weights.

All participating countries used the same general procedure for calculating the survey weights. However, each country developed the survey weights according to its particular probability sample design.

In general, two types of weights were calculated by each country, population weights that are required for the production of population estimates, and jackknife replicate weights that are used to derive the corresponding standard errors.

Population weights

For each respondent record the population weight was created by first calculating the theoretical or sample design weight. Then a base sample weight was derived by mathematically adjusting the theoretical weight for non-response. The base weight is the fundamental weight that can be used to produce population estimates. However, in order to ensure that the sample weights were consistent with a country's known population totals (i.e., benchmark totals) for key characteristics, the base sample weights were ratio-adjusted to the benchmark totals.

Table 5.7 provides the benchmark variables for each country and the source of the benchmark population counts.

Jackknife weights

It was recommended that 10 to 30 jackknife replicate weights be developed for use in determining the standard errors of the survey estimates.

Switzerland produced 15 jackknife replicate weights. The remaining countries produced 30 jackknife replicate weights.

TABLE 5.7 Benchmark variables by country

Country	Source of benchmark counts	Benchmark variables
Bermuda	Census 2000	Age, Gender, Education level
Canada	Census Demography Counts, June-2003	Province, Census geographic area (i.e., CMA/CA), Age, Gender
Hungary	2005 and 2006 demographic data from the Hungarian Central Statistical Office (KSH)	Age, Gender, Educational level, Geographic area
Italy	ISTAT Multipurpose Survey 2002	Region, Age, Gender, Education level, Employment status
Netherlands	Municipal Basic Administration (GBA) as collected by the National Statistical Office (CBS) and Experian database	Age, Education, Purchasing power, House property
New Zealand	2006 Census of Populations and Dwellings	Age, Gender, Educational level
Norway	Norwegian Register of Education (2002 version)	Age, Gender, Education level
Nuevo Leon, Mexico	Census of Population and Housing (2000)	Age, Gender, Education level
Switzerland	Swiss Labor Force Survey (SAKE)	Language region, Age, Gender, Education level, Immigrant status
United States	2003 Current Population Survey, March Supplement	Census region, Metropolitan Statistical Area (MSA) status, Age, Gender, Race/ethnicity, Immigrant status

The following table summarizes the sample sizes and response rates for each participating country.

TABLE 5.8 Sample size and response rate summary

Country	Population aged 16 to 65	Initial sample size (16 to 65)	Out-of-scope cases ¹	Number of respondents ² (16 to 65)	Response rate ³ (16 to 65)
					%
Bermuda	43,274	4,049	745	2,696	82
Canada	21,960,683	35,270	4,721	20,059	66
Hungary	6,760,050	9,178	18,356	5,635	66
Italy	38,765,513	16,727	971	6,853	44
Netherlands	10,974,940	12,734	719	5,617	44
New Zealand	2,634,442	28,702	17,565 ⁴	7,131	56
Norway	2,945,838	9,719	16	5,411	56
Nuevo Leon, Mexico	2,382,454	6,000	36	4,786	80
Switzerland	1,161,735	18,282	5,310	5,120	40
United States	184,260,910	7,045	1,846	3,420	66

1. Out-of-scope cases are those that were coded as residents not eligible, unable to locate the dwelling, dwelling under construction, vacant or seasonal dwelling, or duplicate cases.
2. A respondent's data is considered complete for the purposes of the scaling of a country's psychometric assessment data provided that at least the Background Questionnaire variables for age, gender and education have been completed.
3. The response rate is calculated as number of respondents divided by the initial sample size minus the out-of-scope cases.
4. The reason for the relatively large number of out-of-scope cases in New Zealand is that a screening methodology was used to 'oversample' the Māori and Pacific populations. In the screened portions of the sample, only Māori and Pacific people were treated as in scope.

6.0 Survey Procedures and Data Processing

6.1 Introduction

The ALL procedures were guided by the international guidelines for the administration of the ALL survey. Standard instruments, sampling, collection and processing methodology (including standardized code for Occupation, Industry and Education) are each important components in making the ALL part of an Internationally comparative study program. The following section outlines these procedures and details any deviations from the protocol in Canada. The section will also look at some of the details regarding the post-collection data processing leading up to the creation of the Public Use Microdata File.

The ALL gathered descriptive and proficiency information from sampled respondents through a background questionnaire and a series of assessment booklets containing prose, document, numeracy and problem solving tasks. Survey respondents spent approximately 30 minutes answering a common set of background questions concerning their demographic characteristics, educational experiences, labor market experiences, and literacy related activities. Responses to these background questions make it possible to summarize the survey results using an array of descriptive variables, and also increase the accuracy of the proficiency estimates for various subpopulations. Background information was collected by trained interviewers.

After answering the background questions, the remainder of respondents' time was spent completing a booklet of designed to measure their prose, document, numeracy and problem solving skills. Most of these tasks were open-ended; that is, they required respondents to provide a written answer.

To achieve good content coverage of each of four skill domains, the number of tasks in the assessment had to be quite large. Yet, the time burden for each respondent also needed to be kept within an acceptable range. To accommodate these two conflicting requirements—in other words, to reduce respondents' time burden without sacrificing good representation of the content domain—each respondent was administered only a fraction of the pool of tasks, using a variant of matrix sampling.

6.2 Model procedures manuals and instruments

Each ALL country was given a set of administration manuals and survey instruments to use as a model. Countries were permitted to adapt these models to their own national data collection systems, but they were required to retain a number of key features. First, respondents were to complete the core and main test booklets alone, in their homes, without help from another person or from a calculator. Second, respondents were not to be given monetary incentives for participating. Third, despite the prohibition on monetary incentives, interviewers were provided with procedures to maximize the number of complete background questionnaires, and were to

use a common set of coding specifications to deal with non-response. This last requirement is critical. Because non-completion of the core and main tasks booklets is correlated with ability, background information about non-respondents is needed in order to impute cognitive data for these persons.

6.2.1 Background questions

The model background questionnaires given to each country contained two sets of questions: mandatory questions, which all countries were required to include; and optional questions, which were recommended but not required. Countries were not required to field literal translations of the mandatory questions, but were asked to respect the conceptual intent of each question in adapting it for use. Countries were permitted to add questions to their background questionnaires if the additional burden on respondents would not reduce response rates.

Where the answers to these questions do not compromise the confidentiality of our respondents, the ALL PUMF includes as much of the collected details as possible. Chapter 8 will examine the issues surrounding confidentiality in more details.

6.2.2 Tasks Items

Like the IALS before it, the ALL study is based on the premise that the difficulty of various literacy tasks is determined by certain factors, which are stable across language and culture. Accordingly, all of the ALL countries were given graphic files containing the pool of psychometric items and were instructed to modify each item by translating and adapting the English text to their own language without altering the graphic representation or task characteristics. In many cases, the original item was itself translated into this English model providing everyone with the same starting point. This consistency in the base materials minimized the effects of translation and adaptation errors.

Certain rules governed the item modification process. For instance, some items required respondents to perform a task that was facilitated by the use of keywords. In some cases, the keywords were identical in the question and the body of the item; in others, the keyword was similar but not exactly the same; and in still other cases, the keyword was a synonym of the word used in the body of the item. In another case, respondents were asked to choose among multiple keywords in the body of the item, only one of which was correct. Countries were required to preserve these conceptual associations during the translation process.

Particular conventions used in the items—for example, currency units, date formats, and decimal delimiters—were adapted as appropriate for each country.

To ensure that the adaptation process did not compromise the psychometric integrity of the items, each country's test booklets were carefully reviewed for errors of adaptation.

6.2.3 Standardized non-response coding

It was crucial that the ALL countries managed non-respondent cases in a uniform manner so as to limit the level of non-response bias in the resulting survey estimates.

In ALL, a respondent had to complete the background questionnaire, pass the core block of literacy tasks, and attempt at least five tasks per literacy scale in order for researchers to be

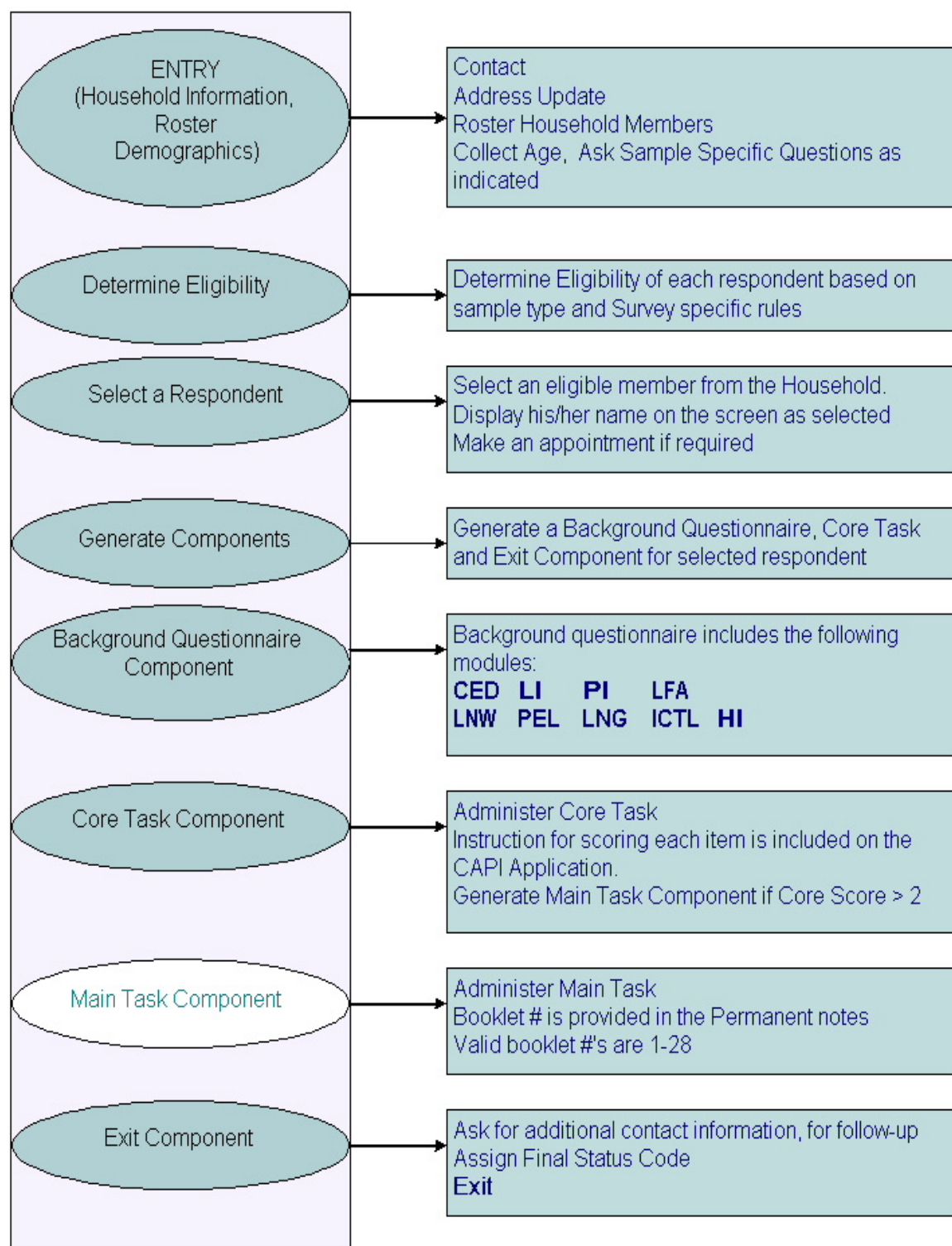
able to estimate his or her literacy skills directly. Literacy proficiency data were imputed for individuals who failed or refused to perform the core literacy tasks and for those who passed the core block but did not attempt at least five tasks per literacy scale. Because the model used to impute literacy estimates for non-respondents relies on a full set of responses to the background questions, IALS countries were instructed to obtain at least a background questionnaire from sampled individuals. They were also given a detailed non-response classification to use in the survey.

Each country was responsible for hiring its own interviewing staff. Thus, the number of interviewers, their pay rates, and the length of the survey period varied among the countries according to their norms and budgets. Each country was provided with a booklet to be used in training interviewers.

In Canada, the ALL was collected by experienced Statistics Canada interviewers using Computer Assisted personal Interviewing technology.

The diagram below graphically depicts the design of the International Adult Literacy and Skills Survey

International Adult Literacy and Skill Survey Main Collection Design



6.3 Scoring

Respondents' literacy proficiencies were estimated based on their performance on the cognitive tasks administered in the assessment. Unlike multiple-choice questions, which are commonly used in large-scale surveys and which offer a fixed number of answer choices, open-ended items such as those used in the ALL elicit a large variety of responses. Because raw data is seldom useful by itself, responses must be grouped in some way in order to summarize the performance results. As they were scored, responses to the ALL open-ended items were ultimately classified as correct, incorrect, or omitted.

The models employed to estimate ability and difficulty are predicated on the assumption that the scoring rubrics developed for the assessment were applied in a consistent fashion within and between countries. Several steps were taken to ensure that this assumption was met. Two of these main steps were the intra-country and inter-country rescores described in the following sections.

6.3.1 Intra-country rescoring

A variable sampling ratio procedure was set up to monitor scoring accuracy. At the beginning of scoring, almost all responses were rescored to identify inaccurate scorers and to detect unique or difficult responses that were not covered in the scoring manual. After a satisfactory level of accuracy was achieved, the rescoring ratio was dropped to a maintenance level to monitor the accuracy of all scorers. Average agreements were calculated across all items. To ensure that the first and second scores were truly independent, certain precautions had to be taken. For example, scorers had to be different persons, and the second scorer could not be able to see the scores given by the first scorer. Scorers who received identical training within a country are expected to be more consistent amongst themselves than with scorers in other countries. Most of the rescoring reliabilities were above 97 per cent. It is important to note that the results were well within the statistical tolerances set for the ALL study and considerably better than those realized in other large-scale studies using open-ended items.

Since intra-country rescoring was used as a tool to improve data quality, score updates were not made to the database. In other words, the agreement data presented here indicate the minimum agreement achieved in scoring. After intra-country reliabilities were calculated, a few scorers were found to be unreliable. These scorers either received additional training or were released. Where scores and rescores differed, the first scores were replaced with correct scores if the inaccuracy was due to a systematic error on the part of the first scorer. In some cases, the scoring guide was found to be ambiguous. In such cases, the scoring guide was revised and the first scores were changed to reflect the revisions, but the second scores were not altered. The second scores were never replaced, even if they were subsequently found to be erroneous.

In sum, the first scores reflect changes and corrections resulting from lessons learned in the intra-country rescoring analysis. The first scores are therefore more accurate and consistent than the second scores, which retain errors and thereby underestimate the rescore reliabilities somewhat. The extent to which the reliabilities are underestimated must be very small, however, given that most of the reliabilities are above 97 per cent. These values indicate that very consistent scoring was achieved by all the participating countries.

6.3.2 Inter-country rescoring

Even after ensuring that all scorers were scoring consistently, fixing ambiguities in the scoring guides, and correcting any systematic scoring errors, it was still necessary to examine the comparability of scores across countries. Accurate and consistent scoring within a country does not necessarily imply that all countries are applying the scoring guides in the same manner. Scoring bias may be introduced if one country scores a certain response differently from the other countries. The inter-country rescoring described in this section were undertaken to ensure scoring comparability across countries.

As noted earlier, responses to the ALL assessment items were scored by each country separately. To determine inter-country scoring reliabilities for each item, the responses of a subset of examinees were scored by two separate groups. Usually, these scoring groups were from different countries. For example, a sample of test booklets was scored by two groups who scored Canada/English booklets and United States booklets. Inter-country score reliabilities were calculated by Statistics Canada, then evaluated by ETS. Based on the evaluation, every country was required to introduce a few minor changes in scoring procedures. In some cases, ambiguous instructions in the scoring manual were found to be causing erroneous interpretations and therefore lower reliabilities.

Using the inter-country score reliabilities, researchers can identify poorly constructed items, ambiguous scoring criteria, erroneous translations of items or scoring criteria, erroneous printing of items or scoring criteria, scorer inaccuracies, and, most important, situations in which one country consistently scored differently from another. In the latter circumstance, scorers in one country may consistently rate a certain response as being correct while those in another country score the same response as incorrect. This type of score asymmetry must be eliminated before the IRT scaling is performed. ETS and Statistics Canada identified such items, while the country in which the scoring problem occurred investigated the plausible causes for such systematic bias in scores. Where a systematic error was identified in a particular country, the original scores for that item were corrected for the entire sample.

Table 6.2 summarizes the inter-country rescore reliabilities of the ALL study before corrections.

Table 6.2 Inter-Country Rescore Reliability Results

Scoring – per cent reliability by domain				
Country pairing (rescoring country – original country)	Psychometric domain			Total
	Prose and document	Numeracy	Problem solving	
	%	%	%	
Canada English – Canada French	95	95	92	95
Canada French – Canada English	95	97	94	95
Norway – Canada	91	93	91	92
Canada – United States	94	97	...	95
United States – Canada	95	97	...	95
United States – Bermuda	91	94	...	90
Bermuda – United States	93	95	...	93
Canada French – Switzerland	95	98	97	96
Switzerland – Canada French	94	96	94	95
Switzerland – Italy	96	98	96	96
Italy – Switzerland	93	97	93	94
Canada – Bermuda	83	83
Canada – Nuevo Leon, Mexico	91	95	...	92
Hungary	94	96	93	94
Netherlands	91	93	93	92
New Zealand	96	97	94	96

... not applicable

1 Quantitative literacy.

6.4 Data capture, data processing and coding

As a condition of their participation in the ALL, countries were required to capture and process their files using procedures that ensured logical consistency and acceptable levels of data capture error. Specifically, countries were advised to conduct complete verification of the captured scores (i.e. enter each record twice) in order to minimize error rates. Because the process of accurately capturing the test scores is essential to high data quality, 100 per cent keystroke validation was needed. Each country was also responsible for coding industry, occupation, and education using standard international coding schemes (International S). Further, coding schemes were provided for open-ended items, and countries were given specific instructions about the coding of such items so that coding error could be contained to acceptable levels.

In order to facilitate comparability in data analysis, each ALL country was required to map its national dataset into a highly structured, standardised record layout. In addition to specifying the position, format and length of each field, the international record layout included a description of each variable and indicated the categories and codes to be provided for that variable. Upon receiving a country's file, Statistics Canada performed a series of range checks to ensure compliance to the prescribed format. Statistics Canada additionally ran consistency and flow edits on the data files received. When anomalies were detected, countries were notified of the problems and were asked to submit cleaned files.

6.5 Derived Variables

A number of derived variables were created to aid research into the antecedents and outcomes of skills. These are described in detail in this section.

Simple Derived Variables

Household information and income

Have dependent children under age 16 living in the household

DV name:	KIDSHOME
DV label:	Have dependent children under age 16 living in the household
DV value labels:	0 'No children under 16 years old living in the household'; 1 'Children under 16 years old living in the household'
Source question(s):	K1: Including yourself, how many people live in your household? K2: Do you have any dependent children living with you in your household? (Children for whom you are financially and/or have sole or joint custody). K3: What is the age of the youngest child in your household?
DV pseudo logic:	K1, K2, and K3 determine whether there is/are dependent child/children under age of 16 living in the household.

7.0 Guidelines for Tabulation and Analysis

This section of the documentation outlines the guidelines to be adhered to by users tabulating, analysing, publishing or otherwise releasing any data derived from the survey microdata tapes. With the aid of these guidelines, users of microdata should be able to produce the same figures as 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.

7.1 Sample Weighting Guidelines for Tabulation

The ALL surveys are based upon complex sample designs, with stratification, 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. In order for survey estimates and analyses to be free from bias, the survey weights must 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 differ 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 poor. Programs for calculating standard errors for simple estimates such as totals, proportions and ratios (for qualitative variables) are provided in the following section.

7.2 Definitions of Types of Estimates: Categorical vs. Quantitative

Before discussing how the ALL data can be tabulated and analyzed, it is useful to describe the two main types of point estimates of population characteristics, which can be generated from the microdata file for the ALL.

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 Albertans at literacy Level 1 on the prose scale or the proportion of Canadians at literacy Level 4 in numeracy are examples of such estimates. An estimate of the number of persons possessing a certain characteristic may also be referred to as an estimate of an aggregate.

Examples of Categorical Questions:

Q: Do you ever watch television or videos in a language other than French or English?

R: **Yes / No**

Q: How would you rate your reading skills in English needed in daily life?

R: **Excellent / Good / Moderate / Poor**

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. They also specifically involve estimates of the form X/Y where X is an estimate of surveyed population quantity total and Y is an estimate of the number of persons in the surveyed population contributing to that total quantity.

An example of a quantitative estimate is the average number of employers that working Canadians had in the past 12 months. The numerator is an estimate of the total number of employers that working Canadians had in the past 12 months, and its denominator is the number of Canadians reporting that they worked in the past 12 months.

Examples of Quantitative Questions :

Q: How many different employers have you had in the past 12 months?

R: **[] [] employer(s)**

Q: How many hours per week did you usually work at this job?

R: **[] [] hours**

7.2.1 Tabulation of categorical estimates

Estimates of the number of people within a given country with a certain characteristic can be obtained from the microdata file by summing the final weights of all records possessing the characteristic(s) of interest.

Proportions and ratios of the form X/Y for a country are obtained by:

- 1) summing the final weights of records having the characteristic of interest for the numerator (X),
- 2) summing the final weights of records having the characteristic of interest for the denominator (Y), then
- 3) dividing the numerator estimate by the denominator estimate.

7.2.2 Tabulation of Quantitative Estimates

Estimates of quantities can be obtained from the microdata file by multiplying the value of the

variable of interest by the final weight for each record, then summing this quantity over all records of interest. For example, to obtain an estimate for a particular country of the total number of different employers that people working part time have had in the past 12 months, multiply the value reported in the question D4 (number of employers) by the final weight for the record, then sum this value over all records with D5=2 (part time).

To obtain a weighted average of the form X/Y , the numerator (X) is calculated as for a quantitative estimate and the denominator (Y) is calculated as for a categorical estimate. For example, to estimate the average number of employers in the past 12 months of people working part time, in a given country

- a) estimate the total number of employers as described above,
- b) estimate the number of people in this category by summing the final weights of all records with QD5=2, then
- c) divide estimate a) by estimate b).

7.3 Skill Level Estimates

The ALL design is an adaptation of a three parameter (PL) Item Response Theory model. The first parameter (A) is the ability of the item to discriminate (sensitivity to proficiency) and the second (B) is its difficulty. A third parameter (C) is the lower asymptote parameter which reflects the possibly non-zero chance of a correct response independent of ability. However, since the ALL test did not generally use any multiple choice type questions, this (C) parameter was fixed at zero throughout, thus transforming the equation into what can now be called a 2PL model. Once the parameters have been calculated, each item can be assigned a Response Probability value of 80 (RP80) which measures the proficiency level needed for a respondent to answer the task with an 80% probability of success.

As noted previously, a respondent's proficiency in the three scales was summarized through the use of the item parameters and the respondent's ability in accordance with the IRT scaling models. The application differed from the norm in that the ALL called for administering relatively few items to each respondent in order to track population levels of proficiency more efficiently. Because the data are not intended to estimate individual levels of proficiency, however, more complicated analyses are required.

Plausible values methodology was used to estimate key population features consistently and to approximate others no less accurately than standard IRT procedures would. In essence, this added dimension requires that the estimation of proficiency be based on a series of five plausible values for each of the three literacy domains. These five plausible values—prose1 through prose5 for the prose scale, doc1 through doc5 for the document scale and num1 through num5 for the numeracy scale and health1 through health5 for the health Literacy scale—have been recoded into plausible levels with values from 1 through 5 reflecting the empirically determined progression of information-processing skills and strategies required to perform increasingly complex tasks. Level 1 is equivalent to scores in the range 0 to 226 (inclusive); Level 2 is equal to scores of 226.0001 through 276; Level 3 goes from 276.0001 to 326; Level 4 includes scores ranging from 326.0001 to 376 and, Level 5 is equivalent to scores greater or equal to 376.0001. For the prose scale, the variables are called plev1 through plev5, for the document scale, these are dlev1 through dlev5 and for the numeracy scale, nlev1 through nlev5.

Due to a difference in the framework, the Problem Solving scale was treated slightly differently. First, the same five plausible values ranging from 0-500 were created (prob1 through prob5), but the level definition was slightly different. For instance, this scale only defines 4 levels of proficiency with level 1 being the weakest and level 4 the highest level of proficiency. Thus, while it is necessary to collapse levels 4 and 5 in order to replicate the published estimates (there are typically too few respondents at level 5 to produce reliable estimates) for the prose, document, numeracy and health scales, this step is not required for the Problem solving domain.

For simple point estimates in either of the five skill domains, it is often sufficient to use the population weight along with one of the corresponding five plausible values (chosen at random).

However, a more precise point estimate can be obtained by taking the average of the five estimates produced from each of the five plausible values, which can be computed as follows:

$T_i = (\sum_j T_{ij}) / 5$, where T_{ij} is a vector of five weighted estimates from each of the five plausible values.

Note that taking an average of the five plausible values, will only produce a valid point estimate, not a valid variance estimate. **All five** plausible values as well as the 30 replicate weights must be used in order to correctly compute design-based variance estimates. Design-based variance estimates are discussed further in section 8.1.2. (Using Plausible Values and Replicate Weights in Calculating Sampling Errors).

7.4 Rounding Guidelines

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

- a) Estimates in the main body of a statistical table are to be rounded to the nearest hundred units using the normal rounding technique. In normal rounding, if the first or only digit to be dropped is 0 to 4, the last digit to be retained is not changed. If the first or only digit to be dropped is 5 to 9, the last digit to be retained is raised by one. For example, in normal rounding to the nearest 100, if the last two digits are between 00 and 49, they are changed to 00 and the preceding digit (the hundreds digit) is left unchanged. If the last digits are between 50 and 99 they are changed to 00 and the preceding digit is incremented by 1.
- b) Marginal sub-totals and totals in statistical tables are to be derived from their corresponding unrounded components and then are to be rounded themselves to the nearest 100 units using normal rounding.
- c) Averages, proportions, rates and percentages are to be computed from unrounded components (i.e. numerators and/or denominators) and then are to be rounded themselves to one decimal using normal rounding. In normal rounding to a single digit, if the final or only digit to be dropped is 0 to 4, the last digit to be retained is not changed.

If the first or only digit to be dropped is 5 to 9, the last digit to be retained is increased by 1.

- d) Sums and differences of aggregates (or ratios) are to be derived from their corresponding unrounded components and then are to be rounded themselves to the nearest 100 units (or the nearest one decimal) using normal rounding.
- e) In instances where, due to technical or other limitations, a rounding technique other than normal rounding is used resulting in estimates to be published or otherwise released which differ from corresponding estimates published by Statistics Canada, users are urged to note the reason for such differences in the publication or release document(s).
- f) Under no circumstances are unrounded estimates to be published or otherwise released by users. Unrounded estimates imply greater precision than actually exists.

8.0 Data Quality

The data quality from any survey can be evaluated by looking at two types of survey errors: sampling error and non-sampling error.

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

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

8.1 Sampling Errors

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

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

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

For example, suppose that, based upon the survey results, one estimates that 16.6% of Canadians are at literacy Level 1 with regard to prose, and this estimate is found to have standard error of 0.013. Then the coefficient of variation of the estimate is calculated as:

$$\left(\frac{.013}{.166} \right) \times 100\% = 7.8\%$$

8.1.1 CV Release Guidelines

One criterion that can be used to determine whether survey estimates are publishable is the coefficient of variation (CV). The CV is the standard error of an estimate expressed as a percentage of that estimate.

Before releasing and/or publishing any estimate from the IALS, users should first determine the quality level of the estimate. The quality levels are acceptable, marginal and unacceptable. Data quality is affected by both sampling and non-sampling errors. However for release purposes, the quality level of an estimate will be determined only on the basis of sampling error as reflected by the coefficient of variation as shown in table 8.1. Nonetheless users should be sure to read section 8 to be more fully aware of the quality characteristics of these data.

First, the number of respondents who contribute to the calculation of the estimate should be determined. If this number is less than 30, the weighted estimate should be considered to be of unacceptable quality. For weighted estimates based on sample sizes of 30 or more, users should determine the coefficient of variation of the estimate and follow the guidelines below. These quality level guidelines should be applied to weighted rounded estimates. All estimates can be considered releasable. However, those of marginal or unacceptable quality level must be accompanied by a warning to caution subsequent users.

Table 8.1 Quality Level Guidelines

Quality level of estimate	Guidelines
1. Acceptable	Estimates have: a sample size of 30 or more, and low coefficients of variation in the range 0.0% to 16.5%. No warning is required.
2. Marginal	Estimates have: a sample size of 30 or more, and high coefficients of variation in the range 16.6% to 33.3%. Estimates should be flagged with the letter M (or some similar identifier). They should be accompanied by a warning to caution subsequent users about the high levels of error associated with the estimates.
3. Unacceptable	Estimates have: a sample size of less than 30, or very high coefficients of variation in excess of 33.3%. Statistics Canada recommends not to release estimates of unacceptable quality. However, if the user chooses to do so then estimates should be flagged with the letter U (or some similar identifier) and the following warning should accompany the estimates: “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.”
--	---

8.1.2 Using Plausible Values and Replicate Weights to calculating Sampling Error

The following section has been liberally copied from the documentation that accompanies the STATTOOL (SAS and SPSS) programs designed by Statistics Canada to help users manipulate the ALL/ALL data. The programs and tools discussed in this section are included with the ALL Public Use Microdata file CD-ROM under the directory called "STATTOOL". While some details of the following section may be more particular to international comparisons of the type facilitated by the ALL Public Use Microdata file, the discussion that ensues will shed some light on the proper usages and practical limits of the ALL data as well.

Calculating point estimates

In this section, we will see how to use the sampling weights (POPWT) to obtain population estimates such as percentages (totals) and means (Calculation of standard errors will be presented in section 6).

All examples will be based on a fictional population with the following characteristics:

Taking the easy way out for preliminary analysis

There is little doubt that the ALL dataset is difficult to manipulate. The 5 Plausible values for the 5 domains (if you include health literacy) along with the 30 replicate weights make the procedures for accurate assessment of standard errors a convoluted affair.

In many instances, simplification of the process, particularly at the exploratory stage would greatly cut down on the processing time required to output the skill estimate analysis.

For this reason, it is recommended that preliminary research use only one of the Plausible values, rather than all five. This is much more accurate than averaging the five plausible values, since it allows for the weighted population distribution to accurately reflect the point estimate. The average of the PV would mask the testing error and, as the population under investigation gets smaller, the estimates will increasingly diverge from the true population distribution.

Of course, once the research is ready for publication, the replicate weights and 5 plausible values should be used to produce the final estimates with accurate standard errors. A full description of this procedure can be found beginning in section 8.1.2.3

Type	Gender	Population Distribution	Sample	
			Distribution (unweighted)	Distribution (weighted)
Rural	Male	40%	30%	38%
	Female	60%	70%	62%
	Total	20%	50%	19%
Urban	Male	51%	45%	50%
	Female	49%	55%	50%
	Total	80%	50%	81%
Total	Male	48.8%	37.5%	47.7%
	Female	51.2%	62.5%	52.3%

	Total	100%	100%	100%
--	-------	------	------	------

From this table, it seems that the male participation was lower than the female participation in both rural and urban areas. Even though nearly 49% of the population is made of males, we have only 37.5% males in the sample. This trend can be observed in both areas. It seems also that the rural area was over allocated with 50% of the sample coming from that area compared to only 20% in the population.

However, once the sampling weights are used, the percentages are quite comparable. How are they calculated?

Percentages (Totals)

The weighted percentage of the males living in rural areas was calculated as follow:

$$\text{Weighted \%} = \frac{\sum_{i=1}^{\text{rural,male}} \text{POPWT}_i}{\sum_{i=1}^{\text{rural}} \text{POPWT}_i} = 38\% \quad \text{where } i \text{ identifies individual } i. \text{ The numerator is an estimate}$$

of the total population of males living in rural areas while the denominator is an estimate of the total population living in the rural areas.

The unweighted percentage was calculated as follow:

$$\text{Unweighted \%} = \frac{\sum_{i=1}^{\text{rural,male}} 1_i}{\sum_{i=1}^{\text{rural}} 1_i} = \frac{n_{\text{rural,male}}}{n_{\text{rural}}} = 30\% \quad \text{where } n_{\text{rural,male}} \text{ is the total number of males living in}$$

rural areas found in the sample and n_{rural} is the total number of people living in rural areas found in the sample.

In the latter, each sampled individual accounts for one while in the weighted version, each sampled unit was given a weight in order to properly and proportionally represent the subgroups in the sample (note that the weighted percentage is a ratio of estimated weighted totals).

Means

For this fictional example, let's say that we also have the average score based on variable PROSE1 as illustrated by the following table:

Type	Gender	Population Distribution	Sample			
			Distribution (unweighted)	Distribution (weighted)	Avg. Prose1 (unweighted)	Avg. Prose1 (weighted)
Rural	Male	40%	30%	38%	260	260.1
	Female	60%	70%	62%	290	289.8
	Total	20%	50%	19%	281.0	278.5
Urban	Male	51%	45%	50%	320	319.7
	Female	49%	55%	50%	330	330.1
	Total	80%	50%	81%	325.5	324.9
Total	Male	48.8%	37.5%	47.7%	296.0	310.7
	Female	51.2%	62.5%	52.3%	307.6	321.0
	Total	100%	100%	100%	303.3	316.1

Here again we see that the weighted means are quite close to the unweighted means as long as one controls by area type. This is not true for the last 3 lines of the table. Let's try to see why. The weighted mean of the males living in rural areas was calculated as follow:

$$\text{Weighted mean} = \frac{\sum_{i=1}^{rural_male} POPWT_i * PROSE1_i}{\sum_{i=1}^{rural_male} POPWT_i} = 260.1 \quad \text{where } i \text{ identifies individual } i. \text{ The}$$

numerator is an estimate of the total score for all males living in rural areas while the denominator is an estimate of the total male population living in rural areas.

The unweighted mean was calculated as follow:

$$\text{Unweighted mean} = \frac{\sum_{i=1}^{rural_male} PROSE1_i}{\sum_{i=1}^{rural_male} 1_i} = \frac{\sum_{i=1}^{rural_male} PROSE1_i}{n_{rural_male}} = 260$$

The unweighted and weighted results will be similar whenever values of PROSE1 don't vary much from one individual to the other and/or values of POPWT behave the same way. This statement doesn't hold for the last three lines of the table. The weighted mean for male is obtained by solving the following equation:

$$\begin{aligned}
 \text{Weighted mean} &= \frac{\sum_{i=1}^{\text{male}} \text{POPWT}_i * \text{PROSE1}_i}{\sum_{i=1}^{\text{male}} \text{POPWT}_i} \\
 &= \frac{(38\% * 19\% * 260.1) + (50\% * 81\% * 319.7)}{47.7\%} = 310.7
 \end{aligned}$$

While the unweighted mean is given by:

$$\begin{aligned}
 \text{Unweighted mean} &= \frac{\sum_{i=1}^{\text{male}} \text{PROSE1}_i}{\sum_{i=1}^{\text{male}} 1_i} \\
 &= \frac{(30\% * 50\% * 260) + (45\% * 50\% * 320)}{37.5\%} = 296.0
 \end{aligned}$$

In the latter, each sampled individual accounts for one while in the weighted version, each sampled unit was given a weight in order to properly and proportionally represent the subgroups in the sample. For example, the 30%X50%=15% of males living in rural areas found in the sample was adjusted by the weights to account for 38%X19%=7.22% of the entire sample which is a much better reflect of what is found in the whole population (Note that the true population proportion of males living in rural areas is 40%X20%=8%).

In conclusion, any statistics computed from sample data should always be done using the sampling weights.

Alternative Sampling Weights

As we saw earlier, the sum of the sampling weights under POPWT within a sample provides an estimate of the size of the population. Although this is a commonly used sampling weight, it sometimes adds to a very large number, and to different numbers from country to country. This is not always desirable. For example, if you want to compute a weighted estimate of the mean achievement in the population across all countries (or sub-populations within a country), using the variable POPWT as your weight variable will lead each country to contribute proportionally to its population size, with the larger countries counting more than small countries. In general, POPWT is not the weight of choice for cross-country analyses. Another consequence of using POPWT is the tendency to inflate results in significance tests when computer softwares are unable to deal correctly with weighted data. We will now see two possible versions of individual sampling weights that address these issues in particular. These versions take advantage on the fact that the same population estimates for means and proportions is obtained whenever you use a weight variable proportional to the population weight (POPWT).

Sum to Constant Sampling Weight (CONSTWT)

It is possible to modify the population weight POPWT such that all countries would contribute the same in a cross-country mean or proportion. This is given by:

$$CONSTWT_{g,i} = POPWT_{g,i} * \left[\frac{100}{\sum_{i=1}^g POPWT_i} \right]$$

for each individual in the group of interest g . The transformation of the weights will be different within each country, but in the end the sum of the variable CONSTWT within each country will be 100. The variable CONSTWT, within each country, is proportional to POPWT multiplied by the ratio of 100 divided by the sum of weights over all individuals in the group of interest. These weights can be used when international estimates are sought and you want to have each country contribute the same amount to the international estimate, regardless of the size of the group of interest in the country (see table below).

Group of interest	Country	Population Count (rural)	Population Estimates	
			Mean PROSE1 (POPWT)	Mean PROSE1 (CONSTWT)
Rural	A	3 700 000	290	290
	B	37 000 000	260	260
	C	7 000 000	300	300
Overall			268	283

Sum to Sample Size Sampling Weight (SMPLWT)

It is possible to modify the population weight POPWT when you want the actual sample size to be used in performing significance tests (within each country). This is given by:

$$SMPLWT_{g,i} = POPWT_{g,i} * \left[\frac{n_g}{\sum_{i=1}^g POPWT_i} \right]$$

for each individual in the group of interest g where n_g is the actual sample size in group g . The transformation of the weights will be different within each country, but in the end the sum of the variable CONSTWT within each country will add up to the sample size in group g . The variable SMPLWT, within each country, is proportional to POPWT multiplied by the ratio of the sample size (n_g) divided by the sum of weights over all individuals in the group of interest. Although some statistical computer software packages allow you to use the sample size as the divisor in the computation of standard errors, others will use the sum of the weights, and this results in severely deflated standard errors for the statistics if POPWT is used as the weighting variable. When performing analyses using such software, it is recommended to use a weighting variable such as SMPLWT as the weight variable. Because of the clustering effect in most country samples, it may also be desirable to apply a correction factor such as a design effect to the SMPLWT variable.

Using the plausible values to compute point estimates

To achieve its goal of broad coverage of the literacy purposes and processes, the ALL /assessment included a range of items arranged into assessment booklets. Each individual participating in the assessment completed one booklet keeping individual response burden to a minimum. ALL used a matrix-sampling design to assign assessment booklets to individuals so that a comprehensive picture of the literacy achievement in each country could be assembled from the components completed by each individual. ALL relied on Item Response Theory (IRT) scaling to combine the individual responses to provide accurate estimates of literacy achievement in the population in each country. The ALL IRT scaling also uses multiple imputation or “plausible values” methodology to obtain proficiency scores in literacy for all individuals, even though each individual responded to only a part of the assessment item pool.

Most cognitive skills testing is concerned with accurately assessing the performance of individual respondents for the purposes of diagnosis, selection or placement. The accuracy of these measurements can be improved by increasing the number of items given to the individual. For the distribution of proficiencies in large population, however, more efficient estimates can be obtained from matrix-sampling design. These designs solicit few responses from each sampled respondent while maintaining a wide range of content representation when responses are aggregated across all respondents. With this approach, however, the advantage of estimating population characteristics is more efficiently offset by the inability to make precise statements about individuals, with the result that aggregations of individual scores can lead to seriously biased estimates of population characteristics.

Plausible values methodology was developed as a way to address this issue by using all available data to estimate directly the characteristics of populations and sub-populations, and then generating multiple imputed scores (called plausible values) from these distributions, which can be used in analyses with standard statistical software. A detailed review of plausible values methodology is given by Mislevy (1991). The main things to retain from this are:

- a) Whenever you want to compute statistics involving scores (like PROSE, DOC, NUMERACY, etc) you don’t have one score value but five score values assigned to each individuals. Each set of plausible values is equally well-designed to estimate population parameters;
- b) These statistics based on scores should always be computed at population or subpopulation levels. They should never be used to do inference at individual level.

Working with Plausible Values

Example1: Estimated median for the variable PROSE in Country A.

For each individual, we don’t have one but 5 scores to deal with as illustrated in the next table

Country A	PROSE1	PROSE2	PROSE3	PROSE4	PROSE5
Individual 1	222	275	300	245	254
Individual 2	289	310	212	250	265
...
Individual n	285	275	243	321	312
Median	285	281	283	279	289

In order to estimate the overall median for PROSE in country A we first have to estimate the median based on the first set of plausible values found under variable PROSE1. We then repeat this first step using PROSE2 through PROSE5 to get a total of five equally good estimates of the median for that country. Since they are all equally good, the next step to do to obtain a single estimate of the median is to average out these five estimates. We get:

$$\text{Overall Median} = (285 + 281 + 283 + 279 + 289) / 5 = 283.4$$

Note that you should not average out scores at the individual level. For example, $(222 + 275 + 300 + 245 + 254) / 5 = 259.2$ is not a good estimate of the variable PROSE for individual 1. This is true for any of the PROSEn variables for a given individual. Score variables should always be interpreted in populations or subpopulations context.

PROSE1 through PROSE5 are the raw ability scores. When these scores are re-grouped into levels, they yield five level variables called for the prose domain PLEV1 through PLEV5 on the PUMF. The variables XPROSE1 through XPROSE5 are a recode of the PLEV1 through PLEV5 variables found on the ALL PUMF such that levels 4 and 5 are collapsed in order to allow sufficient numbers of respondents in each level for accurate analyse. The same can be duplicated for the Document, Numeracy and Health level variables (DLEV1-DLEV5, NLEV1-NLEV5, HLEV1-HLEV5). The Problem Solving levels should be left as found in PSLEV1 through PSLEV5.

Example 2: Estimated logistic regression parameter coefficients (Levels 2 versus level 3 of the dependent variables XPROSE1 to XPROSE5).

From the values found in the previous table this would give:

Country A	XPROSE1	XPROSE2	XPROSE3	XPROSE4	XPROSE5
Individual 1	2	3	3	2	2
Individual 2	3	4/5	2	2	2
...
Individual n	3	2	2	3	3


The first thing we note is that a given individual may be found in different proficiency levels depending on which set of plausible values we are looking at. This does not invalidate the methodology used however. As explained before, in order to get the logistic regression parameter coefficients at the country level, we first have to calculate these parameters based on the first set of plausible values. We then repeat this step 4 times using XPROSE2 through XPROSE5 to get 4 additional sets of estimated parameters as illustrated in the next table:

Country A	XPROSE1	XPROSE2	XPROSE3	XPROSE4	XPROSE5
Intercept	0.124	0.129	0.122	0.125	0.120
Beta 1	1.051	1.059	1.049	1.055	1.060
Beta 2	0.584	0.591	0.545	0.499	0.645
Beta 3	3.222	4.123	3.012	3.542	3.201

Since these sets of estimated parameter coefficients are all equally good, the next step to do to obtain a single set of estimates is to average out these five results. We get:

$$\text{Overall Intercept} = (0.124 + 0.129 + 0.122 + 0.125 + 0.120) / 5 = 0.124$$

$$\text{Overall Beta 1} = (1.051 + 1.059 + 1.049 + 1.055 + 1.060) / 5 = 1.055$$



Overall Beta 2 = $(0.548 + 0.591 + 0.545 + 0.499 + 0.645) / 5 = 0.566$
Overall Beta 3 = $(3.222 + 4.123 + 3.012 + 3.542 + 3.201) / 5 = 3.420$

8.1.3 Estimating Error Variance in ALL

The ALL methodology is a tried and tested method of quantifying skills. As such, a large body of documentation already exists regarding the appropriate ways of estimating variance in such studies. A large part of the text found in this section was borrowed from the IEA PIRLS 2001 User's Guide and adapted to fit the ALL context.

Overview

When analysing data from complex designs such as ALL, it is important to compute correct error variance estimates for the statistics of interest. In ALL this error variance can come from two sources: the sampling process (always present) and the imputation process (whenever the statistics of interest involve proficiency scores). This section describes the methods used to estimate these error variance components.

Estimating Sampling Variance

When data are collected as part of a complex sample survey, analytically there is often no easy way to produce unbiased or design-consistent estimate of variance. A class of techniques called replication methods provides a way to estimate variance for the type of complex sample designs such as those used in ALL.

The basic idea behind replication is to select subsamples repeatedly from the whole sample, calculate the statistic of interest for each subsample, and then use these subsamples or replicate statistics to estimate the variance of the full-sample statistic. Different ways of creating subsamples from the full-sample result in different replication methods. The subsamples are called replicates and the statistics calculated from these replicates are called replicate estimates.

One such method is the jackknife repeated replication (JRR) technique (Wolter, 1985). In ALL, within each country, the full sample was randomly split into 30 subsets of equal or nearly equal size, with each subset resembling the full sample. Replicates are formed by deleting one subset at a time and multiplying the weights for the other subsets by

$$\frac{30}{29}$$

In this manner, 30 replicates are created. This method is also known as the JK1 method. The weights associated with each replicate can be found under variables REPLIC01 through REPLIC30. These weights should only be used to calculate the sampling variance. Point estimates should be calculated as described in the previous section.

Computing Sampling Variance Using the JK1 Method

The basic idea here is to calculate the estimate of interest from the full sample using the weight variable POPWT (or SMPLWT or CONSTWT) as well as each replicate (using the variables REPLIC01 through REPLIC30). The variation between the replicate estimates and the full-

sample estimate is then used to estimate the sampling variance for the full sample. Suppose that $\hat{\theta}$ is the full-sample estimate of some population parameter θ . The sampling variance estimator $\text{var}(\hat{\theta})$ is given by:

$$\text{var}_{\text{smp}}(\hat{\theta}) = \frac{30 \sum_{g=1}^{30} (\hat{\theta}_{(g)} - \hat{\theta})^2}{29} \quad \text{where } \hat{\theta}_{(g)} \text{ is the estimate of } \theta \text{ based on the observations included in the } g\text{-th replicate.}$$

When the statistic of interest involves proficiency scores, it is common practice to base the sampling variance on the first set of plausible values only rather than computing the above expression 5 times and averaging out the 5 estimated sampling variances.

Example 1: Average personal income in country A

The average personal income is given by the following expression:

$$\text{Mean personal income} = \frac{\sum_{i=1}^{\text{Country A}} \text{POPWT}_i * D43_i}{\sum_{i=1}^{\text{Country A}} \text{POPWT}_i} = 27\ 601$$

In order to compute the sampling variance, we have to calculate the following expression 30 times, each one based on the appropriate replicate weight.

$$\text{Mean personal income}_{(1)} = \frac{\sum_{i=1}^{\text{Country A}} \text{REPLIC01}_i * D43_i}{\sum_{i=1}^{\text{Country A}} \text{REPLIC01}_i} = 26\ 983$$

$$\text{Mean personal income}_{(2)} = \frac{\sum_{i=1}^{\text{Country A}} \text{REPLIC02}_i * D43_i}{\sum_{i=1}^{\text{Country A}} \text{REPLIC02}_i} = 26\ 146$$

• • •

$$\text{Mean personal income}_{(30)} = \frac{\sum_{i=1}^{\text{Country A}} \text{REPLIC30}_i * D43_i}{\sum_{i=1}^{\text{Country A}} \text{REPLIC30}_i} = 28\ 965$$

We then simply apply the variance formula given earlier. This gives:

$$\text{var}_{\text{smpl}}(\hat{\theta}) = \frac{29 \sum_{g=1}^{30} (\hat{\theta}_{(g)} - \hat{\theta})^2}{30} = \frac{29}{30} [(26983 - 27601)^2 + (26146 - 27601)^2 + \dots + (28965 - 27601)^2]$$

Finally, the statistic $(\hat{\theta} - \theta) / \text{var}(\hat{\theta})^{1/2}$ is approximately t-distributed with 29 degrees of freedom.

Estimating Imputation Variance

Whenever the statistics of interest involve proficiency scores, there is a need for estimating the imputation variance. As mentioned in previous section, five scores are generated for the same test for individuals participating in ALL. These different scores are referred to as plausible values (PVs). In ALL, each individual was presented with blocks of exercises. The full collection of blocks covers the concepts to be tested, but an individual respondent did not answer questions from all blocks. Using a type of balanced assignment of block of respondents, the full battery of questions was covered when respondents are aggregated. For a group of similar respondents, a Bayesian posterior distribution of scores was estimated. The plausible values for each respondent are realizations from the posterior distribution. These scores are not meaningful for an individual respondent, but when combined can be used to estimate population averages and other population quantities.

Computing Imputation Variance

The general procedure for estimating the imputation variance using plausible values is as follows:

- First estimate the statistic of interest θ , each time using a different set of plausible values (M) and the variable POPWT (or SMPLWT or CONSTWT). Let's call these 5 estimates, $\hat{\theta}_1$ to $\hat{\theta}_5$. The statistic of interest can be anything estimable from the sample data, such as mean, the difference between means, percentiles, regression parameter coefficients, etc.
- Then estimate the overall estimate by averaging out the $\hat{\theta}_m$ where $m=1, 2, \dots, 5$. Let's call this estimate $\hat{\theta}$.
- The imputation variance is computed as:

$$\text{Var}_{\text{imp}}(\hat{\theta}) = \left[1 + \frac{1}{5} \right] \times \sum_{m=1}^5 \frac{(\hat{\theta}_m - \hat{\theta})^2}{4}$$

Estimating the Overall Error Variance

Under ideal circumstances and with unlimited computing resources, the overall error variance would be computed as follows:

$$Var(\hat{\theta}) = \sum_{m=1}^5 \frac{Var_{smp1}(\hat{\theta}_m)}{5} + Var_{imp}(\hat{\theta})$$

Since each of the Var_{smp1} involves calculating the statistics of interest 30 times (each time using a different sampling weight), this shortcut formula can be used instead where sampling variance is estimated from the first set of plausible values only.

$$Var(\hat{\theta}) = Var_{smp1}(\hat{\theta}_1) + Var_{imp}(\hat{\theta})$$

When the statistics of interest do not involve any proficiency scores, the overall error variance formula simply becomes:

$$Var(\hat{\theta}) = Var_{smp1}(\hat{\theta})$$

Degrees of Freedom

$(\hat{\theta} - \theta) / Var(\hat{\theta})^{1/2}$ is approximately t-distributed, with degrees of freedom (Jonhson & Rust, 1993) given by:

$$\nu = \frac{1}{\frac{f_m^2}{4} + \frac{(1-f_m)^2}{29}}$$

where f_m is given by:

$$f_m = \frac{Var_{imp}(\hat{\theta})}{Var(\hat{\theta})}$$

In practice, the number of degrees of freedom is set to 29 (this will work well when f is relatively small, say less than 30%).

Making Comparisons

We will now see how to compute the correct error variance when comparing survey estimates between countries, to the international estimates, and within countries. In order to simplify the text, the estimated mean achievement for the variable PROSE will be considered only. It should be straightforward to generalise this section to any type of survey estimate.

Between Countries

The error variance when comparing the estimated mean achievement for PROSE between country A and B is given by:

$$Var(\hat{\theta}) = Var_{Country A}(\hat{\theta}) + Var_{Country B}(\hat{\theta})$$

For example, say that the estimated mean achievement for country A is 290 with an error variance of 25 and for country B, the estimated mean achievement is 307 with an error variance of 30.25. The difference between these two countries is then 307-290 = 17. The question is : is this difference of 17 points the result of error due to sampling only part of the population combined with the fact that only part of the items were administered? To find the answer to that question, we first have to compute the following statistic (known as the Wald statistic):

$$(\hat{\theta}) / Var(\hat{\theta})^{1/2} = 17 / (25 + 30.25)^{1/2} = 2.287$$

When this value is compared to the critical 95% value from a t distribution with 29 degrees of freedom (2.04), we conclude that there is enough evidence to state that these two countries don't have the same estimated mean achievement.

Note that this approach is also valid when comparing the ALL results to the IALS results.

To the International Estimates

An important published statistics shows your country mean achievement compared to the international mean. The error variance when doing so is given by:

$$Var(\hat{\theta}) = \frac{(N-1)^2 Var_{smp\ A}(\hat{\theta}_A) + \sum_{k=1, k \neq A}^N Var_{smp\ k}(\hat{\theta}_k)}{N^2} + Var_{imp}(\hat{\theta}_A - \hat{\theta}_{int})$$

where N is the number of countries used to compute the international mean, $\hat{\theta}_m$ stands for the estimated mean for country m, and $\hat{\theta}_{int}$ stands for the estimated international mean.

For example, let's consider the following tables:

	Mean Achievement	Sampling Variance
Country A	290	20
Country B	300	22
Country C	286	18
Country D	324	22
International	300	

Mean Achievement	PROSE1	PROSE2	PROSE3	PROSE4	PROSE5
Country A	288	292	292	288	290
International	301	300	299	302	298
$\hat{\theta}_A - \hat{\theta}_{int}$	-13	-8	-7	-14	-8

Here we have that the estimated mean achievement for country A is 290 with a sampling variance of 20 (imputation variance of 5) and the estimated international mean achievement based on 4 countries is 300. The difference between country A result and the international mean is then 290-300 = -10. The question is : is this difference of 10 points the result of error

due to sampling only part of the population combined with the fact that only part of the items were administered? To find the answer to that question, we first have to compute the following statistic:

$$Var(\hat{\theta}) = \frac{(3)^2 20 + (22 + 18 + 22)}{4^2} + \left(1 + \frac{1}{5}\right) \frac{[(13 - 10)^2 + (8 - 10)^2 + (7 - 10)^2 + (14 - 10)^2 + (8 - 10)^2]}{4}$$

$$= 15.125 + 12.6 = 27.725$$

we then compute the Wald statistics:

$$(\hat{\theta}) / Var(\hat{\theta})^{1/2} = 10 / (27.725)^{1/2} = 1.899$$

When this value is compared to the critical 95% value from a t distribution with 29 degrees of freedom (2.04), we conclude that there isn't enough evidence to state that country A is different than the estimated international mean achievement.

Within Countries

Most of the times when comparing subgroups within countries, there is no direct formula for computing the overall error variance like we had in the previous two sections. The main reason for this is that, the samples for the different subgroups are not typically treated as independent for the purpose of statistical tests. Accordingly, a jackknife procedure applicable to correlated samples for estimating the sampling variance of the difference between subgroups should be applied. This involves computing the difference between subgroups once for each of the 30 replicate samples, and five more times, once for each set of plausible values as described earlier (see combining sampling and imputation variances).

However, linear regression models can be easily used to compute these differences. Here's how to compute the difference between men and women for the variable PROSE:

- a) Create dummy variables for the subgroups; let's call MAN the variable that will take value 1 if the respondent is a man and 0 otherwise, and WOMAN the variable that will take value 1 if the respondent is a woman and 0 otherwise. You create as many dummy variables as there are subgroups.
- b) Using POPWT, run a linear regression model with PROSE1 as the dependent variable and either MAN or WOMAN as the independent variable (When there are more k dummy variables with k greater than 2, choose k-1 dummy variables as independent variables). The dummy variable left out will become the reference subgroup.
- c) Using the replicate weights, repeat step b), 30 times.
- d) Using POPWT, repeat step b) 4 times, once for each set of plausible values
- e) Combine information from step b), c), and d) to compute the overall point estimate, the error variance, and the Wald statistics.

For example, say that after creating the dummy variables you get from step b the following result:

Mean PROSE1 = 270 + 18*WOMAN (weighted by POPWT)

This simplifies to Mean PROSE1 = 270 when respondents are males and Mean PROSE1 = 288 when respondents are females. This means that the coefficient in front of the variable WOMAN in the regression model is the difference between women and men while the intercept (270) is the mean achievement for the reference level, men in this case.

From step c you get:

Mean PROSE1 = 271 + 17*WOMAN (weighted by REPLIC01)

Mean PROSE1 = 269 + 19*WOMAN (weighted by REPLIC02)

Mean PROSE1 = 273 + 14*WOMAN (weighted by REPLIC03)

...

Mean PROSE1 = 268 + 21*WOMAN (weighted by REPLIC30)

The sampling variance of the difference between women and men can now be calculated as:

$$\text{var}_{\text{impl}}(\hat{\theta}) = \frac{29 \sum_{g=1}^{30} (\hat{\theta}_{(g)} - \hat{\theta})^2}{30} = \frac{29}{30} [(17-18)^2 + (19-18)^2 + (14-18)^2 + \dots + (21-18)^2] = 19.575$$

From step d you get:

Mean PROSE2 = 271 + 20*WOMAN (weighted by POPWT)

Mean PROSE3 = 269 + 19*WOMAN (weighted by POPWT)

Mean PROSE4 = 273 + 17*WOMAN (weighted by POPWT)

Mean PROSE5 = 268 + 16*WOMAN (weighted by POPWT)

The overall point estimate for the difference can now be computed by averaging out the results over PROSE1 to PROSE5. This gives: $(18+20+19+17+16) / 5 = 18$. The imputation variance is given by:

$$\begin{aligned} \text{Var}_{\text{imp}}(\hat{\theta}) &= \left[1 + \frac{1}{5}\right] \times \sum_{m=1}^5 \frac{(\hat{\theta}_m - \hat{\theta})^2}{4} \\ &= 1.2 \times \frac{[(18-18)^2 + (20-18)^2 + (19-18)^2 + (17-18)^2 + (16-18)^2]}{4} = 3.0 \end{aligned}$$

The Wald statistics becomes

$$(\hat{\theta}) / \text{Var}(\hat{\theta})^{1/2} = 18 / (19.575 + 3.0)^{1/2} = 3.788$$

When this value is compared to the critical 95% value from a t distribution with 29 degrees of freedom (2.04), we conclude that there is enough evidence to state that men mean achievement is different than women mean achievement within country A.

8.1.4 Performing Analyses with the ALL Data Using SPSS

This section presents some basic examples of analyses that can be performed using the sampling weights and scores discussed in previous sections. It also provides details on selected SPSS programs to conduct such analyses, and the results of these analyses. The analyses presented here are simple in nature. The programs compute the percentage of respondents in specified subgroups, the mean achievement for those groups, and the corresponding standard errors (square root of the total error variance) for the percentage and mean statistics.

In our examples, we use macros written in SPSS that can be used to perform any of the analyses that are described in this section. These are general procedures that can be used for many purposes, provided you have some basic knowledge of the SPSS macro language. If you have some programming experience in this statistical package, you will be able to make the necessary modifications to the macros to obtain the desired results.

The SPSS Macros

The four available SPSS macros are described as follows:

JACKMEAN.SPS

This SPSS macro can be used to compute weighted percentages of respondents within defined groups, and their mean (average) on a specified continuous variable. This macro also computes the JRR sampling variances for the percentages and mean estimates. The variable can be any continuous variable in the file.

JACKMEANPV.SPS

This macro can be used in SPSS to compute weighted percentages of respondents within defined groups, and their mean achievement scores on an achievement scale using plausible values. This macro makes use of the plausible values in computing the mean achievement scores. This macro also computes the jackknife repeated replication (JRR) sampling variances for the percentages of respondents within specified groups, and the JRR and imputation variances for the mean achievement scores. This macro should only be used when multiple plausible values are used in the analyses.

JACKREG.SPS

This macro can be used in SPSS to compute linear regression coefficients and their corresponding standard errors within defined groups. This macro can be used with any variable in the analysis but it does not make use of plausible values.

JACKREGPV.SPS

This SPSS macro can be used to compute linear regression coefficients and their corresponding standard errors when using plausible values as the dependent variables within defined groups.

Means and Percentages when Plausible Values are not involved

This section presents example SPSS code that can be used to compute the standard errors for means and percentages of variables other than plausible values. This code is provided in the form of an SPSS program called **JACKMEAN.SPS** that computes the percentages of respondents within subgroups defined by a set of classification variables, the standard errors of these percentages, the means for the groups on a variable of choice, and the standard errors of these means. The standard errors computed by this SPSS macro are taking into account the ALL sample design.

When using this macro, you need to specify a set of classification variables, one analysis variable, the number of replicate weights (if this number is the same for different countries, you can merge the country data sets together, otherwise run the analysis country by country), the replicate weights and the population weight that is to be used for the analysis. You will also need to specify the data file that contains the data to be processed.

You need to know some basic SPSS syntax in order to use the macro effectively. First it needs to be included in the program file where it is going to be used. If you are operating in batch mode, then the macro needs to be called in every batch. If you are using SPSS interactively, the macro needs to be called once at the beginning of the session and it will remain active throughout the session. If the session is terminated or restarted at a later time the macro needs to be called once again. Once the macro is included in a specific session, the word "JACKMEAN" should not be used within that program because doing so will call the macro.

This macro has several parameters. These are:

INFILE The name of the data file that contains the variables necessary for the analysis (If the path location is included as part of the file name, the name of the file has to be enclosed in quotes). Include only the cases that are of interest in the analysis (e.g., respondents with missing variables have to be excluded prior to calling the macro).

CVAR The lists the variables that are to be used to classify the respondents in the data file. This can be a single variable, or a list of variables. It is recommended to always include the variable that identifies the country. At least one variable had to be specified (e.g., CNTRID).

DVAR This is the variable for which means are to be computed. Only one variable can be listed here.

NJKZ This indicates the number of replicate weights that were generated in the data file. When you are working with the data for only one country, you should set the NJKZ argument to as many replicates as are needed in the country (when more than one country data set, make sure all data sets have the same number of replicates).

RPWT The replicate weights in the data files, generally REPLIC01 to REPLIC30. The replicate weights need to be specified in the form "REPLIC01 TO REPLIC30".

WGT The sampling weight to be used in the analysis, generally POPWT.

The simplest way to call the macro is by using the conventional SPSS notation for invoking macros. This involves listing the macro name followed by the corresponding list of arguments for the analysis, each separated by a slash. For example, if the macro is called using the following code:

```
Include "c:\jackmean.sps".
```

```
Jackmean
```

```
      Infile = temp           /
      Cvar  = cntrid          /
      Dvar  = d43              /
      Njkz  = 30               /
      Rpwt  = replic01 to replic30 /
```


Wgt = popwt.

It will compute the mean of personal income (D43) and its standard error, within each country, using the variable POPWT as the sampling weight. The data will be read from the system file TEMP.

The file that contains these results is called FINAL and is saved to the default directory being used by SPSS. The variables that are contained in this file are:

Classification Variables

Each of the classification variables is kept in the resulting file. There is one unique occurrence for each specific combination of the classification variable categories.

Weight Variable

It contains the estimate of the population size of the groups defined by each specific combination of the classification variable categories. In the above example, this variable is called POPWT.

MNX

Contains the means of the variable DVAR for the groups defined by the corresponding combinations of classification variable categories.

MNX_SE

Contains the standard errors of the MNX values computed using the jackknife method.

PCT

Contains the percentages of people in the groups for the classification variable listed last, within the specific combination of the categories defined by the groups initially. In our example, we would obtain the percentages of respondents by country.

PCT_SE

Contains the standard errors of PCT computed using the jackknife method.

The file resulting from using this macro can then be printed using the SPSS procedure of choice. An example is given below.

```
get file = "x:\ALL\ALLDATA.sav"  
  / keep = cntrid gendaa2 d43 popwt replic01 to replic30 .  
  
select if (gendaa2=1 or gendaa2=2) and not(missing(d43)).  
  
save outfile = respondent.  
  
include "c:\ALL\jackmean.sps".  
  
jackmean infile= respondent /
```

```

cvar = cntrid gendaa2      /
dvar = d43                 /
njcz = 30                  /
rpwt = replic01 to replic30 /
wgt = popwt.

```

```
print formats   cntrid gendaa2 n (F6.0) popwt (f10.0) mnx mnx_se pct pct_se (f8.2).
```

```
report format=list automatic / var = cntrid gendaa2 n popwt mnx mnx_se pct pct_se.
```

CNTRID	Gender of Respondent	N	POPWT	MNX	MNX_SE	PCT	PCT_SE
22	1	170	109777	70949.93	7977.68	60.39	3.03
22	2	190	72007	45423.49	6084.25	39.61	3.03

Means and Percentages when Plausible Values are Involved

This chapter presents example SPSS code that can be used to compute the standard errors for mean plausible values and percentages. This code is provided in the form of an SPSS macro called **JACKMEANPV.SPS** that computes the percentages of respondents within subgroups defined by a set of classification variables, the standard errors of these percentages, the means for the groups on one of the achievement scales using plausible values, and the standard errors of these means. The standard errors computed by this SPSS macro are taking into account the ALL sample design and the imputation variance components.

When using this macro, you need to specify a set of classification variables, the name of the plausible values and how many there are, the number of replicate weights (if this number is the same for different countries, you can merge the country data sets together, otherwise run the analysis country by country), the replicate weights and the population weight that is to be used for the analysis. You will also need to specify the data file that contains the data to be processed.

You need to know some basic SPSS syntax in order to use the macro effectively. First it needs to be included in the program file where it is going to be used. If you are operating in batch mode, then the macro needs to be called in every batch. If you are using SPSS interactively, the macro needs to be called once at the beginning of the session and it will remain active throughout the session. If the session is terminated or restarted at a later time the macro needs to be called once again. Once the macro is included in a specific session, the word "JACKMEANPV" should not be used within that program because doing so will call the macro. This macro has several parameters. These are:

INFILE The name of the data file that contains the variables necessary for the analysis (If the path location is included as part of the file name, the name of the file has to be enclosed in quotes). Include only the cases that are of interest in the analysis (e.g., respondents with missing variables have to be excluded prior to calling the macro).

CVAR The lists the variables that are to be used to classify the respondents in the data file. This can be a single variable, or a list of variables. It is recommended to always include the variable that identifies the country. At least one variable had to be specified (e.g., CNTRID).

PVS These are the plausible values to be used in the analysis. The plausible values need to be specified in the form “Plausible Value 1 TO Plausible Value 5” as in “PROSE1 TO PROSE5”. Although in most cases you will want to use all five plausible values, the program will also work when fewer are specified. You should always use at least two plausible values.

NPV This is the number of plausible values that will be used for the analysis. Generally you will want to use all five plausible values for the analysis although under some circumstances fewer can be used (see PVS above).

NJKZ This indicates the number of replicate weights that were generated in the data file. When you are working with the data for only one country, you should set the NJKZ argument to as many replicates as are needed in the country (when more than one country data set, make sure all data sets have the same number of replicates).

RPWT The replicate weights in the data files, generally REPLIC01 to REPLIC30. The replicate weights need to be specified in the form “REPLIC01 TO REPLIC30”.

WGT The sampling weight to be used in the analysis, generally POPWT.

The simplest way to call the macro is by using the conventional SPSS notation for invoking macros. This involves listing the macro name followed by the corresponding list of arguments for the analysis, each separated by a slash. For example, if the macro is called using the following code:

```
Include “c:\jackmeanpv.sps”.
```


```
Jackmeanpv
```

```
      Infile = temp           /  
      Cvar  = cntrid gendaa2  /  
      PVS   = PROSE1 to PROSE5 /  
      NPV   = 5               /  
      Njkz  = 30              /  
      Rpwt  = replic01 to replic30 /  
      Wgt   = popwt.
```

It will compute the mean prose reading achievement and its standard error for males and females within each country, using five plausible values and the variable POPWT as the sampling weight. It will also compute the percentages of males and females within the country, and their corresponding standard errors. The data will be read from the system file TEMP.

The file that contains these results is called FINAL and is saved to the default directory being used by SPSS. The variables that are contained in this file are:

Classification Variables



Each of the classification variables is kept in the resulting file. There is one unique occurrence for each specific combination of the classification variable categories.

Weight Variable

It contains the estimate of the population size of the groups defined by each specific combination of the classification variable categories. In the above example, this variable is called POPWT.

N

Contains the number of cases in the groups defined by each specific combination of categories for the classification variables.



MNX

Contains the means for the first plausible value for the groups defined by the corresponding combinations of classification variable categories.

MNX_SE

Contains the standard errors of the mean for the first plausible value for the groups computed using the jackknife method. This does not include the imputation error.

PCT

Contains the percentages of people in the groups for the classification variable listed last, within the specific combination of the categories defined by the groups initially. In our example, it is the percentage of males and females within each country.

PCT_SE

Contains the standard errors of PCT computed using the jackknife method.

MNPV

Contains the means of the plausible values for the groups defined by the corresponding combinations of classification variable categories.

MNPV_SE

Contains the standard errors for the mean of the plausible values for the groups computed using the jackknife method. This includes the sampling and the imputation components.

The file resulting from using this macro can then be printed using the SPSS procedure of choice. An example is given below.

```

get file = "x:\ALL\ALLdata.sav"
  / keep = cntrid gendaa2 popwt replic01 to replic30 prose1 to prose5.

select if (gendaa2=1 or gendaa2=2) .

save outfile = respondent.

include "c:\ALL\jackmeanpv.sps".

jackmeanpv      infile= respondent      /
                  cvar = cntrid gendaa2      /
                  pvs = prose1 to prose5      /
                  npv=5                    /
                  njkz = 30                  /
                  rpwt = replic01 to replic30 /
                  wgt = popwt.

print formats cntrid gendaa2 (F2.0) n (F4.0) popwt (f7.0)   mnpv   mnpv_se mnx mnx_se
pct pct_se (f6.2).

report format=list automatic margin(1,255)
  / var = cntrid gendaa2 n popwt mnpv mnpv_se mnx mnx_se pct pct_se.

```

CNTRID	Gender of Respondent	N	POPWT	MNPV	MNPV_SE	MNX	MNX_SE	PCT	PCT_SE
22	1	1605	1179970	230.41	1.11	230.39	.99	49.53	.06
22	2	3196	1202504	226.13	.99	226.64	.91	50.47	.06

Regression Coefficients when Plausible Values are not involved

This chapter presents example SPSS code that can be used to compute linear regression coefficients and their standard errors. This code is provided in the form of an SPSS macro called **JACKREG.SPS** that computes the multiple correlation between the specified dependent and independent variables, as well as the regression coefficients and their standard errors. The standard errors computed by this SPSS macro are taking into account the ALL sample design.

When using this macro, you need to specify a set of classification variables, the dependent and independent variables, the number of replicate weights (if this number is the same for different countries, you can merge the country data sets together, otherwise run the analysis country by country), the replicate weights and the population weight that is to be used for the analysis. You will also need to specify the data file that contains the data to be processed.

You need to know some basic SPSS syntax in order to use the macro effectively. First it needs to be included in the program file where it is going to be used. If you are operating in batch mode, then the macro needs to be called in every batch. If you are using SPSS interactively, the macro needs to be called once at the beginning of the session and it will remain active throughout the session. If the session is terminated or restarted at a later time the macro needs to be called once again. Once the macro is included in a specific session, the word "JACKREG" should not be used within that program because doing so will call the macro.

This macro has several parameters. These are:

INFILE The name of the data file that contains the variables necessary for the analysis (If the path location is included as part of the file name, the name of the file has to be enclosed in quotes). Include only the cases that are of interest in the analysis (e.g., respondents with missing variables have to be excluded prior to calling the macro).

CVAR The lists the variables that are to be used to classify the respondents in the data file. This can be a single variable, or a list of variables. It is recommended to always include the variable that identifies the country. At least one variable had to be specified (e.g., CNTRID).

XVAR This is a list of independent variables, at least one, that under the linear regression model will be used as predictors of the dependent variable specified in DVAR. These independent variables can be continuous or categorical, or any other type of coded variable.

DVAR This is dependent variable that under the regression model is predicted by the variable or variables specified by the XVAR parameter. Only one variable can be listed.

NJKZ This indicates the number of replicate weights that where generated in the data file. When you are working with the data for only one country, you should set the NJKZ argument to as many replicates as are needed in the country (when more than one country data set, make sure all data sets have the same number of replicates).

RPWT The replicate weights in the data files, generally REPLIC01 to REPLIC30. The replicate weights need to be specified in the form “REPLIC01 TO REPLIC30”.

WGT The sampling weight to be used in the analysis, generally POPWT.

The simplest way to call the macro is by using the conventional SPSS notation for invoking macros. This involves listing the macro name followed by the corresponding list of arguments for the analysis, each separated by a slash. For example, if the macro is called using the following code:

Include “c:\jackreg.sps”.

Jackreg

```
Infile = temp           /
Cvar = cntrid           /
Xvar = regsex           /
Dvar = d43              /
Njkz = 30               /
Rpwt = replic01 to replic30 /
Wgt = popwt.
```

It will compute the regression equation for the variable REGSEX as a predictor of the personal income. The data will be read from the system file TEMP.

The file that contains these results is called REG and is saved to the default directory being used by SPSS. The variables that are contained in this file are:

Classification Variables

Each of the classification variables is kept in the resulting file. There is one unique occurrence for each specific combination of the classification variable categories.

Mult_RSQ

The squared multiple correlation coefficient for the model.

SS_Res, SS_Reg, SS_Total

The residual, regression and total sum of squares for the model within each group as defined by the classification variables.

Regression Coefficients and Standard Errors (B## and B##.SE)

These are the regression coefficients for each of the predictor variables in the model and their corresponding jackknifed standard errors. The coefficient zero (B00) is the intercept for the model. The other coefficients receive a sequential number starting with 01. This sequential number corresponds to the order of the variables in the list of variables specified in the parameter XVAR.

The file resulting from using this macro can then be printed using the SPSS procedure of choice. An example is given below.

```
get file = "x:\ALL\ALLdata.sav"
  / keep = cntrid gendaa2 d43 popwt replic01 to replic30.

select if (gendaa2=1 or gendaa2=2) .    compute regsex = gendaa2 - 1.

save outfile = respondent.

include "c:\ALL\jackreg.sps".

jackreg          infile= respondent          /
                  cvar = cntrid              /
                  xvar = regsex              /
                  dvar = d43                 /
                  njkz = 30                  /
                  rpwt = replic01 to replic30 /
                  wgt = popwt.

print formats      cntrid (F2.0) n (F4.0) mult_RSQ (f5.3)
                  SS_Total SS_Reg SS_Res (F10.0) B00 B00.SE B01 B01.SE (f6.2) .

report format=list automatic margin(1,255)
  / var = cntrid n Mult_RSQ SS_Total SS_Reg SS_Res B00 B00.SE B01 B01.SE .
```

CNTRID	N	MULT_RSQ	SS_TOTAL	SS_REG	SS_RES	B00	B00.SE	B01	B01.SE
22	360	.023	1.2E+15	2.8E+13	1.18E+15	70950	7977.7	-25526	9560.6

In this example, the variable REGSEX is created by subtracting one from the variable GENDAA2. As a result, males receive a code of 0 and females receive a code of 1 on this variable. In this particular model the variable REGSEX is used to predict the values of the variable D43 (personal income). The model becomes

Personal Income = 70950(7978) for Males,
 Personal Income = 70950(7978) – 25526(9561) for Females.

The numbers in brackets are the standard errors. This means females have on average a personal income that is \$25 526 less than the one for males and \$9561 is the standard error attached to that estimate.

Regression Coefficients when Plausible Values are Involved

This chapter presents example SPSS code that can be used to compute linear regression coefficients using plausible values as the dependent variable and their standard errors. This code is provided in the form of an SPSS macro called **JACKREGPV.SPS** that computes the average multiple correlation between the specified plausible values and independent variables, as well as the regression coefficients and their standard errors. The standard errors computed by this SPSS macro are taking into account the ALL sample design.

When using this macro, you need to specify a set of classification variables, the dependent and independent variables, the number of replicate weights (if this number is the same for different countries, you can merge the country data sets together, otherwise run the analysis country by country), the replicate weights and the population weight that is to be used for the analysis. You will also need to specify the data file that contains the data to be processed.

You need to know some basic SPSS syntax in order to use the macro effectively. First it needs to be included in the program file where it is going to be used. If you are operating in batch mode, then the macro needs to be called in every batch. If you are using SPSS interactively, the macro needs to be called once at the beginning of the session and it will remain active throughout the session. If the session is terminated or restarted at a later time the macro needs to be called once again. Once the macro is included in a specific session, the word “JACKREGPV” should not be used within that program because doing so will call the macro.

This macro has several parameters. These are:

INFILE The name of the data file that contains the variables necessary for the analysis (If the path location is included as part of the file name, the name of the file has to be enclosed in quotes). Include only the cases that are of interest in the analysis (e.g., respondents with missing variables have to be excluded prior to calling the macro).

CVAR The lists the variables that are to be used to classify the respondents in the data file. This can be a single variable, or a list of variables. It is recommended to always include the variable that identifies the country. At least one variable had to be specified (e.g., CNTRID).

XVAR This is a list of independent variables, at least one, that under the linear regression model will be used as predictors of the dependent variable specified by the plausible values. These independent variables can be continuous or categorical, or any other type of coded variable.

ROOTPV This is the prefix used to identify the plausible values for the achievement scale of interest. For example the root of the prose reading plausible values is "PROSE".

NPV This is the number of plausible values that will be used for the analysis. Generally you will want to use all five plausible values for the analysis although under some circumstances fewer can be used (see PVS above).

NJKZ This indicates the number of replicate weights that were generated in the data file. When you are working with the data for only one country, you should set the NJKZ argument to as many replicates as are needed in the country (when more than one country data set, make sure all data sets have the same number of replicates).

RPWT The replicate weights in the data files, generally REPLIC01 to REPLIC30. The replicate weights need to be specified in the form "REPLIC01 TO REPLIC30".

WGT The sampling weight to be used in the analysis, generally POPWT.

The simplest way to call the macro is by using the conventional SPSS notation for invoking macros. This involves listing the macro name followed by the corresponding list of arguments for the analysis, each separated by a slash. For example, if the macro is called using the following code:

Include "c:\jackregpv.sps".

Jackregpv

Infile	= temp	/
Cvar	= cntrid	/
Xvar	= regsex	/
Rootpv	= Prose	/
NPV	= 5	/
Njkz	= 30	/
Rpwt	= replic01 to replic30	/
Wgt	= popwt.	

It will compute the regression equation for the variable REGSEX as a predictor of the plausible values in reading prose. The data will be read from the system file TEMP.

The file that contains these results is called REG and is saved to the default directory being used by SPSS. The variables that are contained in this file are:



Classification Variables

Each of the classification variables is kept in the resulting file. There is one unique occurrence for each specific combination of the classification variable categories.

Mult_RSQ

The squared multiple correlation coefficient for the model.

SS_Res, SS_Reg, SS_Total

The residual, regression and total sum of squares for the model within each group as defined by the classification variables.

Regression Coefficients and Standard Errors (B## and B##.SE)

These are the regression coefficients for each of the predictor variables in the model and their corresponding jackknifed standard errors (with both sampling and imputation components). The coefficient zero (B00) is the intercept for the model. The other coefficients receive a sequential number starting with 01. This sequential number corresponds to the order of the variables in the list of variables specified in the parameter XVAR.

The file resulting from using this macro can then be printed using the SPSS procedure of choice. An example is given below.

```
get file = "x:\ALL\ALLdata.sav"
  / keep = cntrid gendaa2 popwt replic01 to replic30 prose1 to prose5.

select if (gendaa2=1 or gendaa2=2) . compute regsex = gendaa2 - 1.

save outfile = respondent.

include "c:\ALL\jackregpv.sps".

jackregpv      infile      = respondent                        /
                  cvar      = cntrid                          /
                  xvar      = regsex                           /
                  rootpv    = prose                            /
                  npv       = 5                                /
                  njkz      = 30                                /
                  rpwt      = replic01 to replic30              /
                  wgt       = popwt.

print formats      cntrid (F2.0) n (F4.0) mult_RSQ (f5.3)
                  SS_Total SS_Reg SS_Res (F12.0) B00 B00.SE B01 B01.SE (f6.2).

report format=list automatic margin(1,255)
  / var = cntrid n Mult_RSQ SS_Total SS_Reg SS_Res B00 B00.SE B01 B01.SE .
```

CNTRID	N	MULT_RSQ	SS_TOTAL	SS_REG	SS_RES	B00	B00.SE	B01	B01.SE
22	4801	.002	4395153088	10952856	4384200232	230.41	1.11	-4.27	1.50

In this example, the variable REGSEX is created by subtracting one from the variable GENDAA2. As a result, males receive a code of 0 and females receive a code of 1 on this variable. In this particular model the variable REGSEX is used to predict the values of the plausible values reading PROSE. The model becomes

Prose = 230.41 (1.11)	for Males,
Prose = 230.41 (1.11) – 4.27(1.50)	for Females.

The numbers in brackets are the standard errors. This means females have on average a score in prose that is 4.27 less than the one for males and 1.50 is the standard error attached to that estimate.

8.1.5 Performing Analyses with the ALL Data Using SAS

This section presents some basic examples of analyses that can be performed using the sampling weights and scores discussed in previous sections. It also provides details on a selected SAS program to conduct such analyses, and the results of these analyses. The analyses presented here are simple in nature. The program computes the percentage of respondents in specified subgroups, the mean achievement for those groups, the weighted counts of respondents in specified groups, the estimated percentiles for these groups, the regression and logistic regression coefficients and their corresponding standard errors (square root of the total error variance).

In our examples, we use a macro written in SAS that can be used to perform any of the analyses that are described in this section. These are general procedures that can be used for many purposes, provided you have some basic knowledge of the SAS macro language. If you have some programming experience in this statistical package, you will be able to make the necessary modifications to the macros to obtain the desired results.

The SAS Macro

The only SAS available macro is described as follows: **STATTOOL.SAS**

This macro program in SAS can be used to compute several statistics: means, percentiles, frequencies, counts, differences and regressions (standard linear regression, logistic regression and multinomial regression). These statistics are computed within defined groups taking into account the sampling weights. This macro also computes the JRR standard errors with the sampling and imputation components.

Basic Analyses: Means, Percentages, Counts, Percentiles, Regression Coefficients and their Standard Errors

This chapter presents example SAS code that can be used to compute means, Percentages, Counts, Percentiles, Regression Coefficients and their standard errors for any type of variable (whether a plausible values or not). This code is provided in the form of an SAS macro called STATTOOL.SAS that computes these statistics for respondents within subgroups defined by a set of any classification variable (based on plausible values or not). The standard errors computed by this SAS macro are taking into account both sampling and imputation components.

When using this macro, you need to specify a set of classification variables, one analysis variable, the number of replicate weights (if this number is the same for different countries, you can merge the country data sets together, otherwise run the analysis country by country), the

replicate weights and the population weight that is to be used for the analysis. You will also need to specify the data file that contains the data to be processed.

You need to know some basic SAS syntax in order to use the macro effectively. First it needs to be included in the program file where it is going to be used. If you are operating in batch mode, then the macro needs to be called in every batch. If you are using SAS interactively, the macro needs to be called once at the beginning of the session and it will remain active throughout the session. If the session is terminated or restarted at a later time the macro needs to be called once again. Once the macro is included in a specific session, the string "%STATTOOL" should not be used within that program because doing so will call the macro.

This macro has several parameters. These are:

WGT The sampling weight to be used in the analysis, generally POPWT

RWGT The root of the variables specifying the replicate weights in the data files, generally REPLIC01 to REPLIC30. The replicate weights need to be specified in the form "REPLIC".

NREP This indicates the number of replicate weights that were generated in the data file. When you are working with the data for only one country, you should set the NREP argument to as many replicates as are needed in the country (when more than one country data set, make sure all data sets have the same number of replicates).

NPV This is the number of plausible values that will be used for the analysis. Generally you will want to use all five plausible values for the analysis although under some circumstances fewer can be used.

STUDY Put the name of the study (ALL).

CNTRYNO This is the country identifier.

INFILE The name of the data file that contains the variables necessary for the analysis (If the path location is included as part of the file name, the name of the file has to be enclosed in quotes). Include only the cases that are of interest in the analysis (e.g., respondents with missing variables have to be excluded prior to calling the macro).

METHOD This is the statistics you would like to produce. METHOD = mean computes the means of the variable of interest. You can also specify "crosstabs" for crosstabulations, "perc" for percentiles, "diff" for differences, "popest" for population counts, "reg" for standard linear regression, "logistic" for logistic regression and "multinomial" for multinomial logistic regression.

DVAR This is the variable for which means are to be computed. Only one variable can be listed here. Put the name of the variable or only the root if the variable of interest is derived from a set of plausible values.

DVARPV Indicates whether or not the DVAR variable is derived from a set of plausible values. This parameter takes on value 1 if the variable of interest is derived from a set of plausible values and 0 otherwise.

BYVAR This lists the variables that are to be used to classify the respondents in the data file. This can be a single variable, or a list of variables. This parameter defines the subgroups for which means of DVAR are requested.

BYVARPV Indicates whether or not the BYVAR variable(s) is(are) derived from a set of plausible values. This parameter takes on value 1 if the variable of interest is derived from a set of plausible values and 0 otherwise.

In addition to these parameters, 3 parameters, CRITER1, CRITER2, and CRITER3 can be used. They contain one SAS programming statement.

The simplest way to call the macro is by using the conventional SAS notation for invoking macros. This involves listing the macro name followed by the corresponding list of arguments for the analysis, each separated by a comma. For example, if the macro is called using the following code:

```
%include "c:\ALL\stattool.sas";

%stattool (wgt = popwt,
          rwt = replic,
          nrep = 30,
          npv = 5,
          study = ALL,
          method = mean,
          infile = in,
          dvar = prose,
          dvarpv = 1,
          byvar = gendaa2 age3,
          byvarpv = 0 0);
```


It will compute the mean achievement in reading prose using all five sets of plausible values and its standard error, within each group defined by the combination of gender and age categories, using the variable POPWT as the sampling weight. The data will be read from the system file TEMP.

The file that contains these results is called FINALB and is saved to the default directory being used by SAS. There is also a HTML file called FINALB and this one is saved on the C drive of your computer under directory TEMP. This file can easily be accessed using EXCEL from MICROSOFT. The variables that are contained in this file are:

Classification Variables

Each of the classification variables is kept in the resulting file. There is one unique occurrence for each specific combination of the classification variable categories.

ESTIMATE



Contains the means of the variable DVAR for the groups defined by the corresponding combinations of the classification variable categories.

STANDARD ERROR

Contains the standard errors of the ESTIMATE values computed using the jackknife method, including both sampling and imputation components.

PROB > |T|

Gives the probability that a Student statistics be larger than the absolute value of the observed estimate, within the specific combination of the categories defined by the groups initially.

Two examples are given below.

```
libname in "C:\ALL\data";
data in;set in.ALLdata;run;

%stattool(wgt    =    popwt,
          rwgt    =    replic,
nrep    =    30,
5,
method    =    mean,
in,
          dvarpv    =    1,
          byvar    =    gendaa2 age3,
          byvarpv    =    0 0);

study    =    npv    =
          ALL,
infile    =
dvar    =    prose,
```

Study: ALL : , ,

Estimated Means for prose by domain and gendaa2 age3, Based on 5 sets of Plausible Values and 29 D.F.

Controlling for domain

Obs	Domain	GENDAA2	AGE3	estimate	Standard Error	Prob > T
1	ALL	1	1	236.764	1.81853	0
2	ALL	1	2	231.754	1.87374	0
3	ALL	1	3	197.439	4.25769	0
4	ALL	2	1	236.945	2.05748	0
5	ALL	2	2	225.062	1.08269	0
6	ALL	2	3	192.060	2.77765	0

In this example XPROSE1 to XPROSE5 are used as classification variables (they are the Prose PV1 to PV5 each recoded into levels 1 through 5 with levels 4 and 5 collapsed together. We are estimating the mean personal income by level of reading prose.

```
libname in "C:\ALL\data";

data in;set in.ALLdata;run;

%stattool(wgt      =  popwt,
          rwgt     =  replic,
          nrep     =  30,
          npv      =  5,
          study    =  ALL,
          method   =  mean,
          infile   =  in,
          dvar     =  D43,
          dvarpv   =  0,
          byvar    =  XPROSE,
          byvarpv  =  1,
          criter1  =  if d43 < 99999997 );
```

Study: ALL : if d43 < 99999997 , ,

Estimated Means for d43 by domain and xprose, Based on 5 sets of Plausible Values and 29 D.F.

Controlling for domain

Obs	Domain	xprose	estimate	Standard Error	Prob > T
1	ALL	1	39799.22	6167.01	.000000463
2	ALL	2	66658.33	9378.51	.000000081
3	ALL	3	103063.21	27265.80	.000724651
4	ALL	4	172247.60	60242.53	.007788203

Altering the content of the “method” allows for the production of various other statistics as follows.

perc = produces percentiles,

diff = provides differences,

popest = produces population counts,

reg = generates a standard linear regression,

logistic = for a logistic regression and,

multinomial = for the production of multinomial logistic regression coefficient.

All other parameters remains the same as for the examples illustrated above with standard errors due to sampling and test error calculated appropriately for each measure.

8.2 Non-Sampling errors

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

Despite these efforts, non-sampling error is bound to exist in every survey. The following text outlines the most likely sources of this error and its impact on the ALL survey.

8.2.1 Sampling Frame

The use of the 2001 Census insured that the ALL frame was as inclusive as possible and that any exclusions could be effectively calculated into the overall survey design.

8.2.2 Non-response

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.

Total non-response occurred when the interviewer was either unable to contact the respondent, no member of the household was able to provide the information, or the respondent refused to participate in the survey. The national non-response rate for the ALL was around 34%. As described in Section 5.9.1, non-response weighting adjustments were performed to compensate for total non-response. The weighting adjustments were calculated within weighting classes formed by using frame information in the case of non-respondent households (no screener data), and by using screener information in the case of non-respondent individuals (screener completed but no data for the selected respondent). These adjustments were designed to reduce the non-response bias as much as possible with the data that were available.

Partial non-response to the survey occurred, in most cases, when the respondent did not understand or misinterpreted a question, refused to answer a question, or could not recall the requested information. Generally, the extent of partial non-response was small in the ALL.

8.2.3 Response Error

A number of other potential sources of non-sampling error that are unique to the ALL deserve comment. Firstly, some of the respondents may have found the test portion of the study intimidating and this may have had a negative affect on their performance. Unlike “usual” surveys, the ALL test items have “right” and “wrong” answers. Also, for many respondents this would have been their first exposure to a “test” environment in a considerable number of years.

Further, although interviewers did not enforce a time limit for answering questions, the reality of having someone watching and waiting may have, in fact, imposed an unintentional time pressure. It is recognized, therefore that even though items were chosen to closely reflect everyday tasks, the test responses might not fully reveal the literacy capabilities of respondents due to the testing environment. Further, although the test nature of the study called for respondents to perform the activities completely independently of others, situations in the real world often enable persons to sort through printed materials with family, friends and associates. It could be therefore, that the skills measured by the survey do not reflect the full range of some respondents' abilities in a more natural setting.

8.2.4 Scoring

Another potential source of non-sampling error for the ALL relates to the scoring of the test items, particularly those that were scored on a scale (e.g. items that required respondents to write). Special efforts such as centralizing the scoring and sample verification were made to minimize the extent of scoring errors.

9.0 Record Layouts and Univariate Counts

Please refer to the accompanying document 'ALL_Codebook_E.pdf' for the record layout and univariate counts for the data file.

10.0 Principal Participants in the Project

International Direction and Co-ordination

Mr. T. Scott Murray
International Study Director for ALL, Statistics Canada, Ottawa

Mr. Yvan Clermont
International Study Co-ordinator for ALL, first wave of countries, Statistics Canada, Ottawa

Ms. Sylvie Grenier
International Study Co-ordinator for ALL, second wave of countries, Statistics Canada, Ottawa

Mr. Patrick Werquin
International Study Co-ordinator for ALL, OECD, Paris

International Scoring and Scaling

Mr. Irwin Kirsch
Educational Testing Service, Princeton

Mr. Kentaro Yamamoto
Educational Testing Service, Princeton

Ms. Minh-Wei Wang
Educational Testing Service, Princeton

Ms. Julie Eastland
Educational Testing Service, Princeton

National Study Managers

Bermuda	Mr. Crispin Boney <i>Statistics Department, Government of Bermuda, Hamilton</i>
Canada	Mr. Jean Pignal <i>Statistics Canada, Ottawa</i>
Hungary	Janos Wiedermann, Budapest
Italy	Ms. Vittoria Gallina <i>Istituto Nazionale per la Valutazione del Sistema dell'Istruzione, Frascati</i>
Netherlands	Willem Houtkoop <i>Max Goote Expert Center for educational research, Amsterdam</i>

New Zealand	Paul Satherley <i>Ministry of Education, Wellington</i>
Norway	Mr. Egil Gabrielsen <i>Centre for Reading Research, Stavanger</i>
Nuevo Leon, (Mexico)	Mr. Edmundo Guajardo Garza <i>Ministerio de Educación, Monterrey</i>
Switzerland	Mr. Philippe Hertig <i>Office fédéral de la statistique, Neuchâtel</i> Mr. Philipp Notter <i>University of Zürich, Zürich</i>
United States	Ms. Mariann Lemke <i>National Center for Education Statistics, Washington</i> Mr. Eugene Owen <i>National Center for Education Statistics, Washington</i>

Domain Experts and Contributors

Prose and Document

Mr. Irwin Kirsch
Educational Testing Service, Princeton

Mr. Kentaro Yamamoto
Educational Testing Service, Princeton

Ms. Julie Eastland
Educational Testing Service, Princeton

Mr. Stan Jones
Atlantic Heath Promotion Research Center, Yarmouth

Numeracy

Mr. Iddo Gal
University of Haifa, Haifa

Ms. Mieke van Groenestijn
Utrecht University of Professional Education, Utrecht

Ms. Myrna Manly
El Camino College, Palos Verdes

Ms. Mary Jane Schmitt
TERC, Cambridge

Mr. Dave Tout
Language Australia, Melbourne

Mr. Yvan Clermont
Statistics Canada, Ottawa

Mr. Stan Jones
Atlantic Heath Promotion Research Center, Yarmouth



Domain Experts and Contributors

Problem Solving

Mr. Eckhard Klieme

German Institute for International Educational Research, Frankfurt

Mr. Jean-Paul Reeß

LIFE Research and Consult, Bonn

Ms. Anouk Zabal

LIFE Research and Consult, Bonn

Background Questionnaire

Ms. Lynn Barr-Telford

Statistics Canada, Ottawa

Mr. Stan Jones

Atlantic Heath Promotion Research Center, Yarmouth

Mr. Trevor Williams

WESTAT, Rockville

Survey Team, Analysts and Production Team

Ms. Danielle Baum

Statistics Canada, Ottawa

Mr. Richard Desjardins

Statistics Canada, Ottawa

Ms. Sylvie Grenier

Statistics Canada, Ottawa

Mr. John Leung

Statistics Canada, Ottawa

Ms. Carrie Munroe

Statistics Canada, Ottawa

Mr. Owen Power

Statistics Canada, Ottawa
