# SURVEY METHODOLOGY

———

# TECHNIQUES D'ENQUÊTE

June - 1981 - Juin

VOLUME 7

NUMBER 1 - NUMÉRO 1

# C O N T E N T S

Editorial Policy:

The objective of the Survey Methodology Journal is to provide a forum in a Canadian context for publication of articles on the practical applications of the many aspects of survey methodology. The Survey Methodology Journal will publish articles dealing with all phases of methodological development in surveys, such as, design problems in the context of practical constraints, data collection techniques and their effect on survey results, non-sampling errors, sampling systems development and application, statistical analysis, interpretation, evaluation and interrelationships among all of these survey phases. The emphasis will be on the development strategy and evaluation of specific survey methodologies as applied to actual surveys. All papers will be refereed; however, the authors retain full responsibility for the contents of their papers and opinions expressed are not necessarily those of the Editorial Board or the Department. Copies of papers in either Official Language will be made available upon request.

Politique de la rédaction:

La revue Techniques d'enquête veut donner aux personnes qu'intéressent les aspects pratiques de la conduite d'enquêtes, la possibilité de publier sur ce sujet dans un cadre canadien. Les textes pourront porter sur toutes les phases de l'élaboration de méthodes d'enquête: les problèmes de conception causé par des restrictions pratiques, les techniques de collecte de données et leur incidence sur les résultats, les erreurs d'observation, l'élaboration et l'application de systèmes d'échantillonnage, l'analyse statistique, l'interprétation, l'évaluation et les liens entre les différentes phases d'une enquête. On s'attachera principalement aux techniques d'élaboration et à l'évaluation de certaines méthodologies appliquées aux enquêtes existantes. Tous les articles seront soumis à une critique, mais les auteurs demeurent responsables du contenu de leur texte et les opinions émises dans la revue ne seront pas nécessairement celles du comité de rédaction ni de Statistique Canada. On pourra se procurer sur demande des exemplaires d'un article dans l'une or l'autre langue officielle.

---

Submission of Papers:

The Journal will be issued twice a
year. Authors are invited to submit
their papers, in either of the two
Official Languages, to the Editor,
Dr. M.P. Singh, Census and Household
Survey Methods Division, Statistics
Canada, 6th Floor, Jean Talon
Building, Tunney's Pasture, Ottawa,
Ontario, K1A 0T6. Two copies of
each paper, typed space-and-a-half,
are requested. Authors of articles
for this journal are free to have
their articles published in other
statistical journals.

Présentation de documents
pour publication:

La revue sera publiée deux fois l'an.
Les auteurs désirant faire paraître
un article sont invités à en faire
parvenir le texte au rédacteur en chef,
M. M.P. Singh, Division des méthodes
d'enquêtes ménages et du recensement,
Statistique Canada, 6e étage, Edifice
Jean Talon, Parc Tunney, Ottawa,
Ontario, K1A 0T6. Prière d'envoyer
deux exemplaires, dactylographiés à
interligne et demi. Les auteurs des
articles publiés dans cette revue sont
libre de les faire paraître dans
d'autres revues statistiques.

# SURVEY MAINTENANCE - PHILOSOPHY AND PRACTICE[1]

## F. Mayda and P. Timmons[2]

An aspect of surveys not always given adequate consideration
is maintenance.

The scope and importance of survey maintenance are discussed
and a case is made for a more scientific methodological
approach.  Practical applications to various stages of
surveys are illustrated by examples from the Canadian Labour
Force Survey.

## 1.  INTRODUCTION

Survey maintenance is an indispensable part of any continuing survey;
however, its components are usually treated as separate activities
rather than as an overall program.  This ad hoc approach can result
in gaps in the program, inadequate documentation or dissemination of
results, inefficiency, and lack of funding due to inadequate under-
standing of the problems addressed by maintenance.  The recognition
of survey maintenance as a distinct methodological domain can encourage
the application of a more scientific approach.  One such example is
the cost-benefit approach to controls in surveys of Platek and Singh
[1].  This paper will deal with a philosophy of sample maintenance
and will also present some illustrations of various aspects of its
application in a large scale continuous survey.

For the purposes of this paper, survey maintenance can be considered

---

as the sum total of the activities and programs, both regular and
occasional:

a)  which ensure that the survey design is respected in all
    of the operations of the survey,

b)  which measure the quality of operations and of survey data,

c)  which modify or adapt the survey process to meet changing
    requirements.

The maintenance aspect of the design and conduct of large scale continu-
ing surveys seldom receives sufficient recognition. There can be many
reasons for this. Theoreticians, anxious to break new ground and develop
new, and in some sense better, mathematical approaches, find the concept
of maintenance unglamorous and mundane. Managers, concerned with compet-
ing priorities, budgets and production, often fail to see the relationship
of maintenance to operational productivity and question the need for these
expenditures. The very word "maintenance" has the connotation of "just
getting by" or "avoiding deterioration" and seems to imply "no change"
for many people.

In spite of this lack of recognition, maintenance is a most necessary part
of continuing surveys for many good reasons. The requirements of the
survey may change, the conditions in the population being surveyed or the
sample frame itself may change, there may be changes in policy or budget,
new techniques or equipment may become available. Adapting to these
changes and ensuring that quality and efficiency are not compromised is
a regular part of the maintenance of a survey. Even if such obvious
changes do not occur, such things as the turnover of all levels of staff,
the passage of time since principles and procedures were first learned
and the gradual separation of the designers and developers from the
operations staff can lead to a dilution of experience and the possible
degeneration of quality.

The purpose of this paper is neither to try to glamourize survey maintenance nor to develop some all-encompassing theoretical approach, model or package which can be applied to any survey. Rather, the purpose of this paper is to spotlight the role of survey maintenance, to bring to the fore and emphasize how that role applies to large scale continuing surveys and, by demonstrating its function, to encourage a more scientific and theoretical study to be brought to bear on the subject. In order to do this, it must be realized that survey maintenance is more a philosophy than a procedure. Individual programs must be tailored to the needs of the survey.

Most of the discussions on maintenance and the specific examples, which follow, although drawn from the Canadian Labour Force Survey, are relevant for many large scale continuous surveys.

## 2. THE MAINTENANCE PHILOSOPHY

Whenever large scale continuing surveys are planned and developed, significant effort and resources are devoted to implementing the best features that the available money can buy. This ranges from the original sample frame through data collection procedures to final estimation and data dissemination. Once the survey has become operational, especially in the case of surveys used to gather official government statistics, there is a need continually to ensure and to demonstrate the quality of the data and the efficiency of the survey methods and procedures. This ability to demonstrate the validity of the survey is required to allow for quality certification of data, to withstand criticism, to assist the organization in performing program quality audits and to encourage the development of quality improvement programs.

Perhaps the most important feature of this philosophy is to maintain programs which continue in an organized fashion to question operations, procedures and survey materials in order to verify their adequacy.

The anomaly of a good maintenance program is that the more effective
it is in maintaining high quality in the survey process, the less
recognition it may receive as a necessary program. A simple hypotheti-
cal example could be cited. Suppose in a personal interview survey the
interviewers and supervisory staff all know that there is a continuous
and prescribed program of reinterview. This mere presence of the
reinterview program may result in a better standard of data collection.
The more effective this program is, however, the less dramatic will be
the results of the reinterview. When faced with the requirement to
reduce costs it is very tempting for managers to assume that the inter-
view process is properly conducted, as evidenced by the good reinterview
reports, and therefore to cut back on the reinterview program.

The philosophy of a unified survey maintenance program approach implies
a broad scope. Survey maintenance touches on every facet of a survey
from the initial planning to the final dissemination of the survey data.

Many steps or stages can be identified in the process of a survey accord-
ing to the degree of detail one wishes. For convenience, we will broadly
divide the survey process into the following five stages:

    a)   Survey Planning and Design
    b)   Sample Selection and Control
    c)   Data Collection
    d)   Data Capture and Processing
    e)   Estimation and Dissemination

All of the above can be recognized as common to any large scale survey.
When the survey is continuous the planning and design stages are
frequently replaced by periodic improvements and occasional redesigns
and revisions.

## 3. ASPECTS OF MAINTENANCE

Various maintenance programs, to be discussed in detail later, are operative at each of these five stages of a survey. These programs can be classified as Measuring, Controlling and Adapting.

The distinction as to which aspect of maintenance a particular maintenance program falls under is not essential. This is so because often a quality measure, for example, can be used both for diagnostic purposes and as feed-back to operations. What is important is the recognition of the necessity of these aspects in a maintenance program.

### Measuring

Maintenance programs, classified as Measuring, provide certain indicators of performance at various stages of the survey. The measures may be used as a guide to operational control or by data analysts to improve their insight into the reliability of the data and its suitability for particular purposes.

Measurement programs can be identified according to their use:

regulatory: those which serve to measure the conduct of specific aspects of the survey operations.

diagnostic: those which measure how well the survey functions in relation to the survey output.

metadata: measurements of aspects of the survey data used by analysts and managers to evaluate the data itself.

It is understood, of course, that the same measurement may serve more than one of these purposes.

### Controlling

Maintenance programs used in controlling provide measures of survey performance for comparison against standards to identify aspects

requiring correction.  This implies a feed-back mechanism which will
adjust operations to ensure that these standards are met.

Adapting

These programs are essentially means of coping with change, whether due
to changes of objectives or conditions or to the availability of new
methods or equipment.

Examples of maintenance activity in the Canadian L.F.S. for the purpose
of adapting to changed conditions or requirements:

- Parallel run.  The running of two surveys in parallel, the old
  design and the new, to explain any differences and to link up
  the old and new time series.

- Sample size increase.  Differential increase in sample size to
  improve provincial estimates.

- Stabilization program.  To maintain a stable sample size while
  allowing for natural population growth.

- Sample Update.  Partial redesign of new sample units to account
  for unequal growth.

- Sub-sampling in growth clusters.  Reducing the interviewing burden
  on Field operations while allowing for growth.

Data from Maintenance Programs

Data from maintenance programs, whether in the form of quality measures,
feed-back for remedial action or methods review and evaluation can be
considered as outputs from the survey process.  The following schematic
diagram, although by no means complete, serves to illustrate how the
maintenance program can influence all the stages of the survey process.

SURVEY MAINTENANCE

| SURVEY PROCESS | OUTPUTS |
|---|---|

PLANNING
DESIGN
REDESIGN

METHODOLOGY
MANUAL

RESEARCH
FINDINGS

UPDATES

SAMPLE
SELECTION
& CONTROL

SAMPLING
METHODOLOGY

OPERATIONAL
CONTROLS

CHANGES

PROBLEM
SOLVING

MEASURING
CONTROLLING
ADAPTING

DATA
COLLECTION

OPERATIONS
REPORT

MANUALS

INSTRUCTIONS

FEED BACK

PRODUCTION
STANDARDS

DATA
CAPTURE &
PROCESSING

OPERATIONS
REPORT

MANUALS

INSTRUCTIONS

FEEDBACK

QUALITY
REPORTS

ESTIMATION
DISSEMINATION

QUALITY
REPORTS

META DATA

STATISTICAL
EVALUATION

## 4. MAINTENANCE PROGRAMS AND THE SURVEY PROCESS

In order to control and minimize the impact of errors in large scale
surveys, numerous quality control and evaluation programs are tradi-
tionally used. Most of these are familiar to survey methodologists
and variations of them appear in almost any ongoing survey. The more
obvious programs include tabulation and evaluation of non-response,
undercoverage, cost of enumeration, observation, reinterview,
variance/covariance, error rates and so on.

What we would like to do is demonstrate how survey maintenance impacts
on every stage of the survey process. Examples are drawn from the
Canadian Labour Force Survey.

a.  <u>Planning and Design</u>

Obviously if a totally new survey is being planned and designed,
there can be no maintenance program which affects the exercise
directly. However, maintenance programs play two roles in the
planning and design of a continuous survey. The first is in the
fact that survey designers will draw on their own previous experi-
ence and that of others in the area of survey maintenance in order
to evaluate possible features of the survey design. For example,
the type of frame chosen will depend not only on what is available
but also on what experience the designer has been exposed to in
regard to frame maintenance.

The second way in which maintenance programs influence survey design
is in periodic re-designs or programs of survey updating. For
example, the LFS is normally redesigned every 10 years, shortly after
the Decennial Census. The re-design after the 1971 Census was parti-
cularly extensive, incorporating many changes which were suggested
based on the experience of maintenance on the older survey. Details

of the changes made are elaborated in [2]. It is expected that more changes, based on maintenance experience with the current survey, will be introduced in the 1981 Redesign.

No on-going survey regardless of how adequate its initial design was, can remain without change for an extended period of time withot some deterioration. Populations of study change, concepts and objectives are modified, parameters which govern the sample selection become out of date and new procedures and technologies are developed. To prevent deterioration in the level of reliability of survey output, the survey maintenance function must evaluate these new factors constantly and implement required changes.

A specific example of this concerns the up-to-dateness of the LFS sample frame in large cities (SRU areas). After the 1976 Census it became possible to identify population growth in SRU areas from 1971 to 1976. The effect of this growth, which was not uniform even within individual SRU areas, was that size measures used in unequal probability of selection of sampling units became out of date, resulting in increased sampling variances. The changes made to the survey in the 1971 redesign allowed the development of methods to update the design in the SRU areas [3]. By using the information provided from the population comparisons, a special program was introduced to redesign specified sub-units within the SRU [4]. The impact of the re-specification of size measures can be seen in table 1. The effect of the program is to avoid increases in sampling variability due to highly clustered growth. This is particularly significant for estimates at the Census Metropolitan Area level.

## TABLE 1

Increase in number of Random Groups due to SRU Update

December 1977 to March 1981

| Province | No. of sub-units up-dated | Resulting no. of sub-units | Original no. of groups | New no. of groups |
|---|---|---|---|---|
| NFLD | 4 | 4 | 42 | 60 |
| PEI | 5 | 5 | 90 | 112 |
| NS | 15 | 12 | 102 | 128 |
| NB | 12 | 13 | 144 | 214 |
| QUE | 26 | 32 | 162 | 252 |
| ONT | 46 | 53 | 300 | 402 |
| MN | 9 | 8 | 96 | 154 |
| SASK | 6 | 8 | 108 | 194 |
| ALTA | 23 | 39 | 276 | 518 |
| BC | 25 | 27 | 174 | 250 |
| CANADA | 170 | 204 | 1494 | 2284 |

Note: a sub-unit is a contiguous area stratum within a Self-Representing area comprising a number of Random Groups. A Random Group is a random collection of clusters (usually city blocks). The total number of sub-units in the initial design was 734.

b.  Sample Selection and Control

Numerous activities are involved in the selection and control of a sample. In the case of a continuing survey, these relate to the maintenance of the sample frame and the selection and rotation of sample units at various stages.

In the LFS the second last stage of selection is a small well-defined
area called a cluster. All of the dwellings located in the cluster
are identified and listed in the field. The list is stored on a
computerized data base in Ottawa and the final sample consists of a
systematic sample of dwellings drawn by computer from the clusters.
Comparisons of the expected number of dwellings based on the count
when designing the area, to those actually listed frequently show
significant differences. Most often the differences are due to
construction or demolition of dwellings since the time that the
cluster was first defined. In a number of cases, however, differen-
ces were due to incorrect listing or boundary errors. Such errors
result in under or over sampling. To minimize the possibility of
errors in listing clusters, a special program called "Cluster Yield
Monitoring" was established.

Each month, Regional Offices are asked to identify reasons for sig-
nificant differences between the design count of dwellings and the
actual number listed for all newly introduced clusters. The timing
of the program is such that field or design errors can frequently
be corrected before interviewing in selected dwellings has begun.
The following table illustrates some results of the program.

TABLE 2

Cluster Yield Monitoring Program
Number of Exception Clusters Checked (October 1979-December 1980)

| Type of discrepancy | No. of Clusters |
|---|---|
| Valid differences | 968 |
| No correction necessary or no correction possible | 244 |
| Correctable errors | 62 |
| Not determined | 24 |
| Total Exceptions | 1298 |

This evaluation is based on a total of 11804 clusters entering the
active sample during the period.

Another example of maintenance in the Sample Selection and Control stage is stabilization of the sample size. Because of the self-weighting feature of the LFS design and the fact that the Canadian population continues to grow, the sample size would normally continue to grow at the same rate. As a means of holding down survey costs an automated procedure known as Sample Size Stabilization has been developed to keep the sample size from growing [5]. Each month for a specific rotation group and type of area within a province, the number of dwellings selected is compared to a predetermined base figure. Should the number selected be less than or equal to the base, nothing further is done. However, should the number selected exceed the base then the excess of dwellings is systematically dropped from the set of selections. A compensating weight is calculated and applied to all the non-dropped dwellings.

Fluctuations in sample size, due to sampling variability among clusters and unequal growth rates, introduce slight changes in the actual number of dwellings selected each month. However, due to the stabilization program, the net sample size remains fairly stable. The following table illustrates the net sample reductions per month. It can be seen that although the net decrease due to stabilization varies somewhat from month to month, there is an increasing reduction through time compensating for the natural growth in the sample.

## TABLE 3

### SAMPLE SIZE STABILIZATION: OCTOBER 1979 TO MARCH 1981

| SURVEY DATE | NUMBER OF DWELLINGS DROPPED |
|---|---|
| 1079 | 521 |
| 1179 | 544 |
| 1279 | 519 |
| 0180 | 610 |
| 0280 | 544 |
| 0380 | 598 |
| 0480 | 643 |
| 0580 | 667 |
| 0680 | 548 |
| 0780 | 740 |
| 0880 | 677 |
| 0980 | 693 |
| 1080 | 745 |
| 1180 | 868 |
| 1280 | 755 |
| 0181 | 861 |
| 0281 | 847 |
| 0381 | 897 |

At the current rate this amounts to a direct saving in interviewing costs of around $4,500 per month. There are also additional savings due to reduced processing and hiring and training of additional interviewers.

c. Data Collection

This phase of the survey process encompasses all collection activities and the materials used. In the case of the LFS, data is collected by personal and telephone interviews by a large staff of highly trained interviewers. The forms used by interviewers are preprinted for specific households and often the second and subsequent interview show certain data reproduced from the month before. There are many opportunities in this stage of the survey process for

the effective use of maintenance programs to maintain quality and to improve methods and procedures. For example, tabulation and examination of edit changes can lead to improvements in training, changes in questionnaire design or changes in edit rules depending on the results of such analyses. This operation, in the LFS, is called the Field Edit Module and is maintained on a monthly basis.

Other significant modifications can derive from re-interview, observation and cost monitoring programs. It is essential to make results from such programs visible so that their importance can be recognized in order to ensure their continued support.

A phenomenon in large scale probability surveys is the problem of under-coverage. In the LFS the extent to which the survey underrepresents the population is called slippage. Slippage is the accumulated result of many things such as errors in clusters or cluster lists, missed dwellings, missed persons within dwellings, errors in coding and inaccurate population estimates to which the survey estimates are compared.

The continual monitoring of slippage is part of the maintenance program of the LFS. A significant change in the slippage rate triggers remedial action, for example: special list checks or special interviewer instructions.

Another problem in Data Collection is non-response. Continuous monitoring of response rates has shown a consistent trend toward higher non-response due to higher no one at home and temporary absent category during the summer months. In an effort to improve response a procedure known as "Post Survey Week Follow Up" has been developed [6]. In essence this is a special procedure of contacting, mostly by telephone, as many non-responding dwellings as possible one or two days after the normal survey period. Due to considerations such as timeliness, recall length and cost, the procedure is only used within specific restrictions and essentially during the summer months. On occasion, the procedure is also permitted where there are special situations where non-response is expected to be exceptionally high. An example of the sort of improvement that is obtained is shown in the following, Table 4.

## TABLE 4

### POST SURVEY WEEK FOLLOW-UP OTTAWA R. O.

### JULY 1978

| Type of Non-Response | Number ..... | | | Reduction in Non-response Rate (%) |
|---|---|---|---|---|
| | At end of Survey Week | Followed Up | Successful Follow-Up | |
| T | 160 | 117 | 43 | 1.61 |
| N | 46 | 36 | 11 | 0.47 |
| K | 4 | 2 | 2 | 0.04 |
| TOTAL TOTAL | 210 | 155 | 55 | 2.12 |

T = The household was temporarily absent for the entire week.

N = The occupants could not be contacted after several attempts.

K = Circumstances within the household, e.g. sickness, language problems.

d. **Data Capture and Processing**

The next step of the survey process consists of data capture and processing. In any large scale survey, the survey data are transformed into machine readable form; the data are edited and coded and imputations are made. In the LFS, in order to ensure the accuracy of the data capture, a quality control program using complete and sample verification is maintained. The program is designed to ensure that data entry errors do not exceed 3%. Continued monitoring of the program results in very high levels of data capture accuracy and ensures the efficiency of data entry operators.

Even when high levels of data capture accuracy are maintained, the data are still subject to errors which may have been introduced during the field collection. These would represent enumerator or respondent errors and they are detected in the edit process. In the LFS there is a program to identify all data fields where changes have been made during editing. The program is called the Field Edit Module and, as has been mentioned earlier, is used as a feed-back to interviewers, questionnaire designers and editors.

The FEM does not include all errors that might have been made but only those where the data entered (or omitted) causes an edit failure. Nevertheless the FEM results have been shown to be a very good measure of the relative number of errors made. The error rates generated by the FEM are a sensitive indicator of the quality of interviewers' work and also a measure of the awareness of field staff of survey requirements. Since the implementation of the program, there has been steady improvement in the error rates to a currently stable level. Even this stable level shows some improvement, however, each time efforts are made to emphasize accuracy in completing the questionnaire. For example, each time special training sessions for field staff focus on improving accuracy, there is a corresponding improvement in error rates for several occasions immediately after the sessions.

In addition to the FEM analysis of edit failures, the editing section monitors the number of error-containing records each month. Results of this monitoring are regularly discussed with field staff to keep them aware of the need to minimize these errors. In cases where error rates show abnormal increase, there is feed-back to individual Regional Offices with identification of the specific types of errors and suggestions for eliminating them in future.

Table 5 shows how the maintenance of this program has contributed to reducing the error rates over time for the household record (F03) and individual questionnaire (F05).

TABLE 5

EDIT ERROR RATES BY FORM TYPE

1977 - 1980

| F03 | JAN. | FEB. | MAR. | APR. | MAY | JUNE | JULY | AUG. | SEP. | OCT. | NOV. | DEC |
|------|------|------|------|------|-----|------|------|------|------|------|------|-----|
| 1977 | 5.8 | 4.7 | 4.8 | 5.2 | 5.1 | 5.3 | 5.3 | 5.4 | 5.4 | 5.1 | 4.9 | 4.6 |
| 1978 | 4.8 | 4.2 | 4.0 | 3.4 | 3.6 | 3.3 | 3.1 | 3.0 | 3.1 | 2.9 | 2.8 | 2.9 |
| 1979 | 2.3 | 2.3 | 2.3 | 2.6 | 2.3 | 2.6 | 2.3 | 2.5 | 2.3 | 1.9 | 2.1 | 2.2 |
| 1980 | 1.8 | 1.9 | 1.8 | 1.9 | 1.9 | 1.9 | 2.1 | 1.8 | N/A | 1.9 | 1.8 | 1.9 |

| F05 | JAN. | FEB. | MAR. | APR. | MAY | JUNE | JULY | AUG. | SEP. | OCT. | NOV. | DEC. |
|------|------|------|------|------|-----|------|------|------|------|------|------|------|
| 1977 | 20.3 | 18.9 | 17.7 | 17.7 | 18.0 | 15.6 | 16.7 | 15.6 | 15.4 | 16.4 | 15.2 | 13.6 |
| 1978 | 14.4 | 13.8 | 13.5 | 14.5 | 14.3 | 11.9 | 12.6 | 11.6 | 11.3 | 12.7 | 11.1 | 9.8 |
| 1979 | 10.3 | 9.7 | 9.8 | 7.7 | 9.9 | 9.7 | 9.6 | 8.9 | 8.4 | 7.5 | 8.5 | 8.4 |
| 1980 | 7.7 | 7.4 | 7.1 | 8.6 | 8.0 | 8.0 | 8.1 | 7.1 | N/A | 8.2 | 7.6 | 6.4 |

e.  Estimation and Dissemination

In this final phase of the survey process, the estimates are published and distributed to the various users. In a sense the data leave the hands of the survey methodologists and enter the domain of the data analysts and policy makers.

Survey data, especially those from large scale continuing surveys, are usually collected for two reasons. The first is to document and

chart the phenomenon of interest. Thus, in the LFS, the survey data serve to provide a comprehensive summary statement about the labour force activity of the Canadian population. The second reason is for policy makers to combine data from various sources to evaluate existing social policy, to predict trends and to devise new policy aimed at improving the social situation.

It is not immediately obvious how maintenance programs can affect this phase of the survey process. Estimation procedures are usually fixed in that they depend on the probability design of the survey. Changes are only made if the probabilities of selection change. We have an example in the LFS. As mentioned earlier, in an effort to put a limit to natural sample size growth, a procedure of "stabilization" was implemented. The effect was the requirement to add a special weight to compensate for the sample reduction. This weight was incorporated into the estimation process.

Other changes too are being incorporated which must be considered part of the regular maintenance function. A program is currently under way to expand from two-digit to three-digit occupational codes to respond to requests from users for more detail. Such a change will not be implemented without careful assessment of the impact on editing and coding operations, data processing, tabulation and print-ing and finally estimation of data reliability.

By carefully monitoring the data processing operations, it has been possible to improve head office processing schedules to such a degree that the press release date for LFS data has been advanced. It has moved forward four days from the Tuesday at the start of the third week after survey week to the Friday at the end of the second week.

Perhaps the most neglected aspect of survey maintenance and quality evaluation in general is its impact on the uses to which data is put after publication. Too frequently implicit assumptions are made to the effect that the published data is "true" without considering outside factors which should temper any analysis. Most data

analysts and users recognize the existence and significance of
sampling errors, but may be less appreciative of the uncertainties
in the data caused by non-sampling errors and non-respose (missing
data).

Maintenance programs are effective in reducing not-sampling errors
and non-response (missing data). They also provide valuable inform-
ation (metadata) which should be considered by the analyst in using
the data. Difficulties in data estimation and evaluation caused by
non-sampling errors and non-response are most difficult to deal with.
They must then be controlled by programs of prevention and this in
addition to maintaining operations is the purpose of maintenance
programs.

## 5. CONCLUSION

This has been a rather short and not very detailed overview of some of
the maintenance programs of the LFS. Particular attention has been
given to some of the less well known programs in an effort to demons-
trate how fundamental they can be and how they form a part of overall
survey maintenance.

We hope the foregoing makes the case for the pervasiveness and importance
of survey maintenance and its contribution to better and more useable
statistics. We recognize that there are, however, substantial costs
involved and methods need to be developed to produce dynamic indicators
similar to cost-variance studies used in sample design. A start in this
direction has been made by Platek and Singh [1].

The relative value and cost of the various procedures should be used to
control the scope, incidence and intensity of the components of the survey
maintenance program.

RESUME

La coordination est un aspect des enquêtes auquel on n'accorde pas toujours suffisamment d'attention.

Les auteurs analysent l'envergure et l'importance de la coordination des enquêtes et préconisent une approche méthodologique plus scientifique. Des applications pratiques à diverses étapes des enquêtes sont illustrées par des exemples tirés de l'Enquête sur la population active du Canada.

REFERENCES

[1] Platek, R. and Singh, M.P., (1980), "Cost Benefit Analysis of Controls in Surveys", Symposium on Survey Sampling, Ottawa.

[2] Statistics Canada (1977), "Methodology of the Canadian Labour Force Survey", Catalogue No. 71-526.

[3] Platek, R. and Singh, M.P., (1978), "A Strategy for Updating Continuous Surveys", Metrika, Vol. 25, pp 1-7.

[4] Drew, D., Choudhry, H. and Gray, G., (1978), "Some Methods for Updating Survey Frames and Their Effects on Estimation", Proceedings of the American Statistical Association, pp 62-71.

[5] Drew, D. (1977), "LFS Sample Size Stabilization", Technical Memorandum, Census and Household Survey Methods Division, Statistics Canada.

[6] Gower, A. (1979) "Non-Response in the Candian Labour Force Survey", Survey Methodology, Statistics Canada, Vol. 5, No. 1. .

# IMPUTATION IN SURVEYS : COPING WITH REALITY[1]

## I.G. Sande[2]

In surveys a response may be incomplete or some items may be inconsistent or, as in the case of two-phase sampling, items may be unavailable. In these cases it may be expedient to impute values for the missing items. While imputation is not a particularly good solution to any specific estimation problem, it does permit the production of arbitrary estimates in a consistent way.

The survey statistician may have to cope with a mixture of numerical and categorical items, subject to a variety of constraints. He should evaluate his technique, especially with respect to bias. He should make sure that imputed items are clearly identified and summary reports produced.

A variety of imputation techniques in current use is described and discussed, with particular reference to the practical problems involved.

## 1. INTRODUCTION

Everyone who has been involved in surveys knows that life would be very easy if only the respondent had read the textbook. If he had, he would know that he is allowed to respond correctly and completely, or not to respond at all. He is not allowed to respond incorrectly or incompletely. Unfortunately, the respondent has not read the textbook. Furthermore, if you call him back to correct the data or fill in missing information, he may not be very co-operative. More often than not, the cost of calling back is simply too high to be carried out generally.

---

So reality might look like this:

## TABLE 1

## IMPORTANT CANADIAN SURVEY

| Record No | Identification Classification | Weight | Variables 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|
| 1 | X | $W_1$ | A | a | y | 3.1 | 4.3 |
| 2 | X | $W_2$ | A | a | z | 4.6 | 2.8 |
| 3 | X | $W_3$ | A | b | y | - | 1.1 |
| 4 | X | $W_4$ | B | b | z | 2.3 | 4.6 |
| 5 | X | $W_5$ | B | c | y | 4.9 | 2.3 |
| 6 | X | $W_6$ | B | b | - | 3.2 | 3.6 |
| 7 | X | $W_7$ | C | - | x | 3.0 | - |
| 8 | X | $W_8$ | C | - | y | - | 1.2 |
| 9 | X | $W_9$ | C | a | - | 0.0 | 2.4 |
| 10 | X | $W_{10}$ | - | b | y | - | 1.4 |

Edits: $A \wedge a \Rightarrow$ Not x.

$B \wedge b \Rightarrow$ Not y.

Var 4 + Var 5 $\leq$ 10.

Var 4 $\geq$ 0, Var 5 $\geq$ 0.

This survey has both categorical and numeric items, and there are three constraints (edits) on the items which must be satisfied. We notice that of 10 records, four (1, 2, 4, 5) are complete. If we look hard, we might also notice that the "missings" are informative: a low value of Variable 5 is associated with a missing Variable 4.

Our primary problem is that we have to produce tabulations of population estimates, e.g. Variable 1 x Variable 2 x Classification Variables, or Variable 4 x Classification variables. Although we might be able to write down all the estimates we think we have a need for in our publication, we know that after the publication comes out, we are going to get a large number of requests for tabulations and estimates which we have not anticipated.

How, then, are we to deal with the partial non-response? The possibilities are:

(i) Ignore all the records with missing values. This may result in loss of a great deal of data, since many records may be affected. Furthermore, "missings" are seldom random and the procedure would almost certainly lead to biased estimates.

(ii) Publish "unknowns" as a category. This is a little better than (i); but still ignores the partial information about the missing value which may be available in the other variables. Frequently, the users of the data will make adjustments for the "unknown" categories without being able to look at the microdata and with little knowledge of the data collection process.

(iii) Adjust (reweight) each table or estimate, ignoring the missings in each case. This is a variation of (i) which may give rise to inconsistent tables in the sense that no complete data set corresponds to the set of estimates because of the constraints on the data.

(iv)  Fill in the blanks in each record with plausible and consis-
tent values.  This is called imputation.

To sum up, partial non-response arises in two ways:

(i)   A record (i.e. the total response for a single survey unit)
contains one or more missing values because (after all possible
checking and follow-up) the data are unavailable.

(ii)  A record is inconsistent in the sense that its component items
do not satisfy natural or reasonable constraints (known as
edits) and one or more items are designated unacceptable (and
therefore are artificially "missing").

To cope with the "missing value" problem in an expeditious manner, values
are frequently imputed for the missing items so that the data set is
"completed".

The estimation of individual values in a data set is not a new problem.
It is the direct descendant of the "missing observation" problem in
ANOVA and the "incomplete data" problem in multivariate analysis.  How-
ever, though imputation is not an optimal solution to the "missing
value" problem in surveys when any particular estimates are considered,
it may just be the least bad of the feasible solutions for general
purposes.

## 2.  THE GENERAL IMPUTATION PROBLEM

What are the "facts of life" facing the unwilling imputer?  No matter
what method of imputation he opts for, the following problems must be
dealt with:

(i)  The close relationship between editing and imputation.

    (a)  If a record fails an edit, it is not always obvious which
fields are faulty, but some basis must be established for
deciding which fields to change.  Does one change all the
fields involved in a failed edit?  Some of them may be invol-
ved in other edits which do not fail.  Does one change the
least number of items, as recommended by Fellegi and Holt [9],
or adopt a policy of "least change", whatever that means?
Or does one adopt the "principle of expedience" : deleting
that configuration which makes imputation easy?

These are non-trivial problems.  The mathematical analysis of
edits and the identification of fields to be changed when
several edits have been failed, is a very subtle problem.
Fellegi and Holt did the first systematic work on categori-
cal or coded data and their methods have been implemented at
Statistics Canada and used (with modifications, see [11]) in
the Census of Population.  The parallel work for numerical data
with linear edits has been carried out by Gordon Sande at
Statistics Canada using optimization techniques [20] and the
development of techniques for the combined numerical and cate-
gorical data problem is seen as feasible.

    (b)  When it has been decided which fields must be imputed
(because they are missing or must be changed) it is obvious
that the imputed data must satisfy the edits, i.e. the comp-
leted record must be consistent.  This requirement often
eliminates the mathematically elegant imputation schemes and
reduces the mathematical tractability of the problem to zero.
Since complex edits make the imputation procedure hard, the
theoretical analysis of such procedures is virtually impos-
sible.  Therefore edits are usually ignored in theoretical
work on the properties of imputation techniques.

(ii) The marginal and joint distributions of responses are almost certainly different from those of the underlying population. In the case of numeric data, such distributions are unlikely to be normal. Transformations to normality (or less pronounced skewness) result in transformations of the edits which makes them more difficult to deal with.

(iii) The pattern of missing fields varies from record to record. In an n-field record (excluding the identifiers and classification variables), there are $2^n-1$ possible patterns of fields to impute. Some imputation schemes (I do not know if any have been seriously implemented) seek to specify a separate imputation procedure for each pattern; but if n is large, this idea soon gets out of hand:

(iv) The imputer does not usually have much time to fiddle with the data after they have come in. Most survey data should be processed promptly to be useful and in some cases (such as many at Statistics Canada) the time constraints are severe. Therefore the method of imputation should be precisely specified before the processing begins. Furthermore, the statistician usually has little, if any, test data to work on before the data collection begins. Historic data cannot always be trusted to look like current data in any but the most general respects. For example we may believe that X is proportional to Y on the basis of historic data; but the proportion $\frac{X}{Y}$ may change from year to year. On the other hand, the circumstances governing the joint occurrence or non-occurrence of X and Y may be similar over time, a fact which can be exploited in testing imputation procedures.

(v) Imputation does not solve any specific estimation problem more satisfactorily than classical estimation techniques for incomplete data, and it may do a lot worse. The trouble is that if one can optimally estimate a particular $\underset{\sim}{\theta}$ using some (correct) distributional assumptions and a (correct) model, one hasn't solved the problem for $\underset{\sim}{\phi}$. One has to start again. If one combines $\underset{\sim}{\theta}$ and $\underset{\sim}{\phi}$,

one may have an unwieldy problem. By the time one has optimally estimated all the parameters one can think of, one may have a set of estimates which is not consistent with any possible data set. And then someone may find a $\underset{\sim}{\psi}$ to be estimated. By imputing a consistent value for each missing item one can estimate any of the usual population parameters (means, totals, ratios, differences, proportions, correlations) very easily, although possibly with no guaranteed precision.

(vi)  It is generally hard to know how to estimate the variance of estimates when some data is imputed. If the amount of imputed data is very small, the usual estimates will do. In some circumstances, mathematical or empirical studies in a vaguely related situation may be available.

(vii)  The imputer is faced with ethical problems if the microdata are ever going to be given out. At the very least, he must plan to identify the imputed items on all copies of the data and publish the proportions of imputations in each field as part of a discussion of data quality when the primary results are published. Alternatively, he may choose to give out edited, but unimputed, versions of the data set. In this case, the secondary users may do their own imputations and get results which are inconsistent with each other and the original.

Which data set should be analyzed? The question really is: What do you mean by analysis? If one wants to explore relationships between variables, the use of imputed data could be prejudicial, not to mention misleading. For simple estimation purposes, as we have pointed out, the imputed set reduces the headache. And we could argue that if the data are so bad that the presence of imputed data could influence the analysis significantly, then the data are not worth analyzing.

After considering these problems we may conclude that the imputer needs
a procedure which

  (i)  will impute plausibly and consistently provided only that the
       non-missing data satisfy the edits;

 (ii)  will preserve the underlying distributions in the data or, at
       least, reduce the response bias and preserve the relationships
       between items as far as possible;

(iii)  will work for (almost) any pattern of missing items;

 (iv)  can be set up and tested ahead of time;

  (v)  can be evaluated in terms of data quality and impact on precision
       of the estimates.

Particular techniques of imputation vary in their ability to meet these
requirements.


## 3. METHODS OF IMPUTATION

Planning ahead is to be recommended.  If one can guess the fields most
likely to cause problems, it will pay to pick up a correlated variable
on the questionnaire or from auxiliary sources.  For example, it may be
hard to get information about household income, but easy to get an
estimate of square feet of living space or some other correlate of income.
The store manager may not want to disclose his gross income;  but one can
count the number of cash registers.  How this information is used depends
on the circumstances.

Techniques of imputation vary from naive to sophisticated.

(i)     Use of ad-hoc values.  Each case may be treated differently in
        a manual procedure, or a few rules of thumb are formulated on
        the basis of "experience" and hunches, and often without the
        encumbrance of real facts.  These are used to fill in the blanks.

        For example, in a business survey we may have the rule for imputing
        the value of closing inventory:  if gross income (GI) is less than
        or equal to $25,000, set closing inventory (CI) to 0; if GI is
        greater than $25,000, set CI equal to 5% of GI minus net income
        or 0, whichever is larger.  In many ways this rule appears quite
        reasonable, provided GI and net income are always available,
        especially if the 5% came from last year's survey.  If it is dirty
        it is at least quick and not too damaging if only a small percentage
        of the records are affected.

        Rules of this type can be formulated to force compliance with the
        edits.  They are also compatible with the simplest of data pro-
        cessing systems.  However, they are subjective and may not reflect
        reality.  The effects on the underlying distributions are often
        unpredictable and non-response bias is not necessarily reduced.
        Evaluation may be impossible.

(ii)    Post-stratify and use the post-stratum marginal mean or another
        typical value (e.g. the mode in the case of a categorical vari-
        able), making sure that there are sufficient data in each post-
        stratum.  In the numeric case, this is equivalent to item by item
        reweighting.

        In the closing inventory example of (i), we might post-stratify
        by gross income, net income, industry, region, etc.  If we create
        too fine a grid or too many data are missing, some collapsing may
        be necessary to ensure that there are enough good data in each
        cell (see [8]).

This technique may run into trouble with the edits. If this seems likely, some modification may be in order (such as letting the edits define the post-strata). Like the method of ad-hoc values, it is very simple, if it works; but will create spikes in the marginal distributions and may be biased. However, in the numeric case variance estimates are generally available.

(iii) Model the relationships between the variables. A popular idea has been to use the conditional mean given the items present, modified to account for the information in the incomplete records assuming normality, or some generalization of this idea (e.g. [3], [7], [12]). However, normality is not usually a plausible assumption and it does not take the edit structure into account. I have not seen any theory worked out for non-normal cases and I am not aware of any application to missing survey data except for test purposes (e.g. Huddleston and Hocking in [1], pp. 88-93).

In one survey at Statistics Canada, about 160 items are collected (from administrative documents) for a fairly small sample of businesses and 5 major items are collected from other sources for the entire population. For various reasons (mainly the ease of arbitrary tabulation of estimates) it is desired to impute the 160 items for the non-sampled businesses. A ratio-type imputation is used, after stratification by size and industry:

$$\hat{x}_i = \frac{\sum\limits_{P} x_j}{\sum\limits_{P} Y_j} \ Y_i$$

where $x$ is related to major item $Y$ and the $i$th record requires imputation. $P$ is the sample of complete records with all 160 items present. Because of the structure of the data, the edits are automatically satisfied; but the imputations do not reflect the real structure of the data which have a lot of zero values. In other words, the imputed records are not realistic and the marginal

distributions are distorted.  On the other hand, the principal
estimates (which are just ratio estimates) are quite acceptable
and permit variance estimation.  In this case the ratio-type
imputation is used because it is easy and convenient, not because
it is a good model.  The effort that would go into fitting a model
would be prodigious and one may well never achieve a good fit.

Thus modelling is an elegant solution which will probably reduce
bias.  On the other hand achieving a good fit may require a great
deal of effort or one may have to tolerate a bad fit, and there
may be problems with edits.  Furthermore, one may find that the
assumed model becomes "built into" the data and may be recovered
by other researchers later, unless steps are specifically taken
to prevent this.

(iv)  Use of historic data, such as last month's or last year's response
for the same unit, if available.  This technique is in common use in
monthly surveys where the same units are surveyed in consecutive
months, for variables which are not expected to change often.  Of
course, the assumption is that one did get a response for the
particular item at some stage and when one has carried a value
forward for several months in a row, one perhaps ought to do some
investigation into what is going on.

(v)  Use a proxy data from another source.  This means that another file,
perhaps of administrative data such as medical or tax records, is
available with the unique identifiers required for matching to the
survey file and that this file includes an equivalent item which
can be used as a proxy for the missing survey item (e.g. [10]).

If an exact match is not available (possibly because the identifiers
have been removed for reasons of confidentiality), one may be content
with a statistical match on classification fields such as age, sex,

and place of birth.  For example, one may use last year's sample
survey as a source of data for statistical matching and imputation
for this year's survey.

Most statistical matching is used for linking different data files
to extend data sets (see e.g. Radner in [1], pp. 108-113).  The
idea of statistical matching is closely related to the hot deck
and nearest neighbour techniques discussed in (vi) and (vii)
below.

(vi)  Use of the current survey data as a source of matched individual
      data records from which one (the donor) is selected at random to
      supply values for missing items in a particular deficient record.
      Procedures of this type are often called hot deck procedures; but
      there is no agreement on the definition of hot deck procedures in
      the literature.  I will take it to mean an imputation procedure
      which uses records from the current survey to supply missing
      values and involves a random or pseudo-random choice.  There seem
      to be two main variants currently in use, both directed mainly at
      categorical data:

      (a)  The sequential hot deck, used in the U.S., for example, in the
           Current Population Survey and the Census of Population.  Here
           the data are processed one record at a time.  To impute a
           field or group of fields A, a cross-classification (matrix)
           of several other related fields (B,C,D...) is defined.  For
           each cell in this classification, that value of A is retained
           which occurred in the last record processed with the corres-
           ponding values of B,C,D.... .  Thus, as the file is processed,
           the values in the individual cells of the B,C,D... matrix
           change.  When a record lacking a value for A occurs, it receives
           the value currently in the cell of the matrix which matches
           its own values of B,C,D...  If two such records (missing A,
           but with the same values of B,C,D...) occur consecutively,
           the same value of A will be imputed in each case.

The ordering of the file may not be random, so that the record used as a donor is not chosen at random. In fact, it may not be advantageous to randomize the file, thereby exploiting the correlations between nearby records to improve the imputation.

The matching fields (and therefore the imputation matrix) vary with the fields to be imputed, so that many matrices must be maintained. In those cases where imputation of a single field might result in an edit failure after imputation, a set of related fields is deleted and imputed together.

Because different fields are imputed from different imputation matrices, several donors may be involved in completing a single deficient record and this may be a source of some concern.

Each imputation matrix must be initialized, using historic data or ad-hoc values. On the other hand, the imputation can be done at one pass and is not difficult computationally.

(b) The random choice procedure used by the Canadian Census and Labour Force Survey. Here an imputation matrix is not maintained; but the set of records with the required values in the matching fields is identified and the donor is chosen at random from these to supply the missing items to the deficient record.

In the Canadian Census, an attempt is made to impute all missing items on a deficient record using a single donor. If this fails, a field-by-field hot deck is tried, in which several donors may be involved [11].

The choice of matching fields in both sequential and random choice procedures must be made considering likely sources of variation, linkage through edits and the number of complete or eligible records available as potential donors in each cell. If too many fields are used for matching, the number of

-

-

- 34 -

potential donors may be too small; if too few fields are used for matching, there is a risk of a poor match or edit failure in the imputed record.

With hot deck methods, the variance of the estimates in simple cases is known to be larger than the variance of the usual expansion estimates of means and totals (e.g. [2]). However, there may be a reduction in bias.

(vii) Use of the current survey data as a source of individual data records with similar characteristics to supply values for missing items. Unlike the hot deck procedures in (vi), these procedures are appropriate for use with numeric data. I shall call them nearest neighbour procedures rather than hot deck procedures because the value in the matching fields must be similar (not the same) and the element of randomness in the choice of donor may be absent.

The hot deck procedures discussed in (vi) run into trouble when numeric fields are linked by edit constraints and matching must be done on them. Occasionally the problem can be dealt with by splitting the range of the variable, e.g. age, into intervals and coding the intervals; but consider the problem of imputing the age of a child from the age of a parent.

For purely numeric data with linear edits, a prototype system at Statistics Canada locates the m "nearest" complete records to a particular deficient record. An attempt to complete the deficient record using fields from the nearest of the m neighbours is made. If the tentatively completed recipient record passes the edits, the imputation is complete. Otherwise, the next nearest neighbour is tried, and so on. If none or the m neighbours will do, the imputation fails and further processing is required [20].

In this type of imputation, the use of suitable data transformations can make the imputation proceed more smoothly. It also helps to insert additional edits so that extreme observations are not admitted as donors (special arrangements can be made for them).

The method requires an efficient search algorithm; but the choice of distance function is not crucial and one which is simple computationally is advisable.

It is possible that particular records will be used as donors much more often than others. Another nearest neighbour type of imputation system developed at Statistics Canada for the imputation of mixed numeric and categorical data, incorporates the number of times a particular record has been used as a donor into the distance function, so that the distance increases with the number of previous donations [5].

Nearest neighbour procedures can be converted into hot deck procedures by choosing the donor record at random from m nearest neighbours instead of taking the nearest satisfactory record. Both types of procedure can be regarded as a form of non-parametric regression.

With numeric matching, the variance would be hard to calculate since the match is deterministic given the data.

(viii) Use of hybrid methods. In fact, to my knowledge, no complex imputation problem is handled by a single imputation procedure. Some ad hoc imputations are usually combined with more sophisticated methods so that the job gets done expediently. Typically, some items are imputed one way and others another way and then some cleaning up is done. In one case [22], the occurrence of zeros in a particular variable was modelled. Those missing cases not imputed as zero through the model were imputed by hot deck.

Various devices may be employed to expedite the imputation. Among these are:

(i) Formulation of the edit procedures to reduce the number of possible missing configurations. More fields than necessary

are deleted, but consistent imputation is easier. For
example, if the edit is $A + B + C \leq X$, failure of the edit
may result in the deletion of all fields A,B,C and X or just
A,B,C rather than only one of these fields. Obviously this
is an option to be used with extreme caution since information
is destroyed.

(ii)   Transformation of the data. It is sometimes more natural to
       impute proportions than absolute numbers and often the edits
       transform neatly to permit this. For the purpose of numerical
       hot decks or nearest neighbour procedures, the distannce function
       function is often better formulated in terms of transformed
       variables than the originals which may be very skew. In terms
       of the original variables, "nearness" in one part of the space
       may be quite different from "nearness" in another.

(iii)  Dividing the record into segments and imputing one segment at
       a time. Each pass is conditional on the preceding ones being
       complete. This makes the imputation task less formidable and,
       in those cases where matching is required, allows different
       appropriate matching procedures to be used at each stage [5].
       A related device is to attempt a global imputation first and,
       where this fails, to try a stage by stage imputation [11].
       If all else fails, we can end with an ad-hoc procedure to tie
       up the loose ends.


## IV  EVALUATION OF IMPUTATION PROCEDURES


In evaluating an imputation procedure, the relevant concerns are bias and
variance of the estimates (means, ratios, etc.) not the ability of an
imputation procedure to guess missing values of individual items correctly.

The theoretical treatment of imputation procedures is generally confined to fairly simple cases, ignoring edit constraints (e.g. Bailar and Bailar in [2] and [15], pp. 422-447; Schaible in [15], pp. 170-187; Platek and Gray, [17]; Szameitat and Zindler, [23]). Empirical work deals either with the comparison of different imputation methods (e.g. [6], [8], [22]); or with the performance of a particular technique under different conditions ([5], [10]). Various edit and imputation strategies are compared by Nordbotten in [13]. Other studies simply attempt to examine the impact of imputation [14], or summarise current practice [18].

Since the scope for theoretical work is limited to fairly simple data and imputation procedures, it seems that, in general, imputation procedures must be evaluated by simulation. This usually means selection or creation of a clean data set (no items missing) to act as a population, the creation of artificial "missings" in biased and unbiased modes and at different rates, and studying the performance of the imputation process over several replicates of each case. The quality (bias, variance), in relationship to the rate and bias of "missings", of the resulting esti- mates may then be assessed. Particular imputation procedures will allow variants of this basic recipe: for example, in a sequential hot deck, replicates may be generated by re-ordering the data set rather than by regenerating a complete set of "missings" as required by nearest neighbour techniques.

Rubin [19] advocates the routine production of several sets of imputed values under different models or sets of assumptions, as part of the regular data processing. This leads to estimates of the "imputation error", that part of the error due to imputation, in the actual data and the effects of different models can be studied. The method which is applicable to only a limited variety of imputation techniques, including hot deck, has been used experimentally.

In general, the estimation of the "imputation error" under normal pro- duction conditions will be very difficult; but it is better to use

approximations obtained from a simulation study than nothing at all.

Whatever the method of imputation, the actual imputation process should be carefully monitored. In the simplest cases this means recording data about the missing items which were subsequently imputed : the number of records in which any imputation is made, the number requiring one (two, three, etc.) item(s) to be imputed, the number of records missing specific variables (or possibly combinations of variables), statistics breaking down the imputations into those due to item non-response and those due to edit failure. For imputations made using a decision tree (the imputation being conditional on other fields and the relationships between them), the number of imputations made in each branch of the tree should be recorded. For a nearest neighbour procedure one also wants to know, for example, how many times each record was used as a donor, which donor was involved in a particular imputation, how many attempts were required to complete a record and what the value of the distance function was. And of course one wants a listing of any records failing to be completed. (It is also equally important to monitor the editing process which precedes the imputation).

## V. CONCLUSION

This is not the first paper on imputation in surveys (e.g. [4], [16], [18], [23]) nor will it be the last. The activity has been going on for a long time under such disguises as "automatic error correction" and used to be considered as part of data processing rather than statistical methodology. Now the survey statisticians are getting involved and the subject is being discussed in the literature and at meetings. Predictably, the open discussion of imputation has dismayed some of the more classical statisticans.

Reality does not consist of the data at the end of the chapter (like the iris data) and normal distributions: it consists of 20,000 long

forms filled out by 20,000 businessmen with other things on their minds, or several million census returns filled out by individuals who want to get back to the newspaper or the TV. These people want to be co-operative; but if the information requested isn't handy or has been forgotten, they omit the question or make up a response, and they also make mistakes. The survey people have to extract as much sense as possible from the results and they try to do a respectable and ethical job.

Reality also consists of the almost unlimited and unpredictable demands which are made on some data sets. These should be satisfied in a consistent way. And reality is the fact that even the simplest survey, properly run, is a complex operation and one does not want to increase the complexity any more than one has to.

I believe that the real problem of imputation is the interaction with editing. Very little of the literature deals with this problem. Szameitat and Zindler [23] and Nordbotten []3] touch on the subject. The "Canadian School" led by Fellegi and Holt ([9]; [5], [11], [20] and even [21]) discuss it (with little empirical work), while, by and large other writers do not, preferring to simplify the problem so that it is amenable to mathematical analysis or empirical study. This does not suggest to me that the effort is wasted, but that the problem of studying the properties of imputation procedures under realistic conditions is a very difficult one. And one must admit that there are some one-question surveys to which the available results might be applicable.

I hope that we will see more empirical work on data sets with complex edit constraints. We need to know much more about how imputation procedures compare with each other and we need guidance about how to optimize the performance of a specific type of procedures. So far, we have only scratched the surface.

## RESUME

Dans les enquêtes, il arrive qu'une réponse soit incomplète
ou que certains éléments soient incompatibles ou encore, que
des éléments puissent manquer, comme dans le cas de
l'echantillonnage à deux phases. Il peut alors être utile
d'imputer des valeurs aux éléments manquants. Même si cette
méthode n'offre pas une solution particulièrement bonne à un
problème d'estimation donné, elle permet cependant la production
d'estimations arbitraires d'une façon cohérente.

Le statisticien enquêteur sera peut-être aux prises avec un
mélange d'éléments numériques et qualitatifes qui seront
assujettis à une variété de contraintes. Il doit évaluer sa
technique, en particulier en ce qui concerne le biais, et
veiller à ce que les éléments imputés soient nettement identi-
fiés et que des rapports sommaires soient produits.

L'auteur décrit diverses techniques d'imputation utilisées à
l'heure actuelle et elle accorde une attention particulière
aux problèmes pratiques en cause.

## REFERENCES

[1]    Aziz, F. and Scheuren,  F., Imputation and Editing of Faulty
       or Missing Data, 1978, U.S. Department of Commerce, Bureau of
       the Census (papers presented at the 1978 meetings of the
       American Statistical Association, almost all appearing in the
       Proceedings of the Section of Survey Research Methods).


[2]    Bailar, J.C. III and Bailar, B.A., "Comparison of Two Procedures
       for Imputing Missing Survey Values". Proceedings of the Section
       on Survey Research Methods, American Statistical Association, 1978,
       pp. 462-467. Also [1], pp. 67-75.


[3]    Beale, L.M.L., and Little, R.J.A., "Missing Values in Multivariate
       Analysis", Journal of the Royal Statistical Society, Series B,
       Vol. 37, 1975, pp. 129-145.

[4]     Chapman, D.W., "A Survey of Nonresponse Imputation Procedures".
        Proceedings of the Social Statistics Section, American Statistical
        Association, 1976, pp. 245-329.


[5]     Colledge, M.L., Johnson, J.H., Pare, R., and Sande, I.G., "Large
        Scale Imputation of Survey Data". Proceedings of the Section on
        Survey Research Methods, American Statistical Association, 1978,
        pp. 431-436. Also, Survey Methodology, Statistics Canada, 1978,
        Vol. 4, No. 2, pp. 203-224; and [1], pp. 102-107.


[6]     Cox, B.C. and Folsom, "An Empirical Investigation of Alternate
        Item Nonresponse Adjustments". Proceedings of the Section on
        Survey Research Methods, American Statistical Association, 1978,
        pp. 219-223; also [1], pp. 51-55.


[7]     Dempster, A.P., Laird, N.M. and Rubin, D.B., "Maximum Likelihood
        from Incomplete Data via the E M Algorithm". Journal of the
        Royal Statistical Society, Series B, Vol. 39, 1977, pp. 1-11.


[8]     Ernst, L.F., "Weighting to Adjust for Partial Nonresponse",
        Proceedings of the Section on Survey Research Methods, American
        Statistical Association, 1978, 1978, pp. 468-472. Also [1],
        pp. 87-91.


[9]     Fellegi, I.P., and Holt, D.A., "Systematic Approach to Automatic
        Edit and Imputation". Journal of the American Statistical
        Association, Vol. 71, 1976, pp. 17-35.


[10]    Ford, B.L., "Missing Data Procedures: A Comparative Study".
        Proceedings of the Social Statistics Section, American Statistical
        Association, 1976, pp. 324-329.

[11]    Hill, C.J., "A Report on the Application of a Systematic Method
        of Automatic Edit and Imputation to the 1976 Canadian Census".
        Proceedings of the Section on Survey Research Methods, American
        Statistical Association, 1978, pp. 474-479.  Also Survey
        Methodology, Statistics Canada, Vol. 4, 1978, pp. 178-202;
        and [1], pp. 82-87.

[12]    Hocking, R.R., and Marx, D.L., "Estimation with Incomplete Data:
        an Improved Computational Method and the Analysis of Mixed Data".
        Communications in Statistics - Theory and Methods, Vol. A8,
        1979, pp. 1155-1182.

[13]    Nordbotten, S., "The Efficiency of Automatic Detection and Correction
        of Errors in Individual Observations as Compared with Other Means
        of Improving the Quality of Statistics".  Bulletin of the
        International Statistical Institute.  Proceedings of the 35th
        Session, Vol. 16, 1965, pp. 417-441.

[14]    Ono, M., and Miller, H.P., "Income Nonresponse in the Current
        Population Survey".  Proceedings of the Social Statistics Section,
        American Statistical Association, 1969, pp. 277-288.

[15]    Panel on Incomplete Data of the Committee on National Statistics/
        National Research Council, Symposium on Incomplete Data:
        Preliminary Proceedings, 1979.  U.S. Department of Health,
        Education and Welfare, Social Security Administration Office of
        Policy, Office of Research and Statistics.

[16]    Platek, R., "Causes of Incomplete Data, Adjustments and Effects".
        Survey Methodology, Statistics Canada, Vol. 6, 1980, pp. 93-132.

[17]    Platek, R., and Gray, G.B., "Nonresponse and Imputation",
        Survey Methodology, Statistics Canada, Vol. 4, 1978, pp. 144-177.

[18]    Pritzker, L., Ogus, J., and Hansen, M.F., "Computer Editing
        Methods - Some Applications and Results". Bulletin of the
        International Statistical Institute. Proceedings of the 35th
        Session, Vol. 16, 1965, pp. 442-465.

[19]    Rubin, D.B., Multiple Imputations in Sample Surveys - A Phenomo-
        logical Bayesian Approach to Nonresponse". Proceedings of the
        Section on Survey Research Methods, American Statistical
        Association, 1978, pp. 20-28.

[20]    Sande, G., "Numerical Edit and Imputation". International
        Association for Statistical Computing, 42nd Session of the
        International Statistical Institute, December 1979.

[21]    Sande, I.G., "A Personal View of Hote Deck Imputation Procedures".
        Survey Methodology, 1979, Statistics Canada, Vol. 5, pp. 238-258.

[22]    Schieber, S.J., "A Comparison of Three Alternative Techniques
        for Allocating Unreported Social Security Income on the Survey
        of the Low-Income Aged and Disabled". Proceedings of the Section
        on Survey Research Methods, American Statistical Association, 1978,
        pp. 212-218. Also, [1], pp. 44-50.

[23]    Szameitat, K., and Zindler, H.J., "The Reduction of Errors in
        Statistics by Automatic Corrections". Bulletin of the
        International Statistical Institute. Proceedings of the 35th
        Session, Vol. 16, pp. 395-417.

# REDESIGNING CONTINUOUS SURVEYS IN A CHANGING ENVIRONMENT

## M.P. Singh and J.D. Drew[1]

Survey organizations undertake periodic redesigns of
continuous surveys.  Reasons for such redesigns related
to changes in information needs to be satisfied by the
survey and changes in public awareness and attitudes
towards surveys are discussed in the context of the
redesign of the Canadian Labour Force Survey following the
1981 Census.  In particular, the importance of close
dialogue between users of the survey data and design
statisticians at the early stages of the redesign process
in order to establish survey objectives is stressed.

## 1.  INTRODUCTION

Data from decennial censuses in addition to serving the need of their
primary users, serve as one of the frequently used tools in design-
ing new surveys and by far the most important tool for redesigning
(designing) large scale continuous surveys of population and housing.
For instance, the Canadian Labour Force Survey  (CLFS), a monthly
survey of 55,000 households across Canada [ 12 ], has been redesigned
following each decennial Census.  Two of the primary reasons for
these post-censal redesigns are to update the sample design to reflect
changes in population characteristics and boundaries of Census units,
and to incorporate improved methodologies  such as sample selection
and estimation procedures.  Also the redesigns provide a unique
opportunity to respond to a) changes in information needs to be

---

[1] M.P. Singh and J.D. Drew, Census and Household Survey Methods Division,
Statistics Canada.

satisfied by the survey, and b) changes in public awareness and atti-
tudes towards surveys and other factors affecting data collection.  In
this paper, the discussions are focused on items a) and b) in the context
of the redesign of the CLFS following the 1981 Census.  They may
be found relevant for other similar surveys.

With regard to information needs, it should be mentioned that at the time
of the revision of the CLFS carried out during the 1970's  [ 10 ], a great
deal of emphasis was placed on more data and increased reliability of data
at the provincial level.  However, new and important uses of CLFS data
have emerged since that time, such as the legislated use of CLFS estimates
in determining the eligibility for benefits under the Unemployment
Insurance Program, administered by Employment and Immigration Canada.
The redesign currently being planned for will represent the first occasion
to consider such new data requirements.  Along with the uses of labour
force data, uses of the Labour Force Survey vehicle itself for obtaining
other socio-economic data have greatly increased in recent years.  After
briefly discussing the process of identifying the survey objectives
in general terms, the discussion in section 2 will cover the following
three specific situations for meeting data demands:

 

    i)    Improved monthly data at sub-provincial levels through reallocation
        of sample within the provinces.

    ii)   Reliable data on quarterly and annual basis at smaller levels of
        aggregations through alternate rotation patterns.

    iii)  Increased demand for socio-economic data vis-à-vis current survey
        capacity.

The third section of the paper discusses the impact of the data collection
method-item (b)- on development of the design of the survey, and
emphasizes the importance of timely decisions on the procedure to be

adopted. The seventies have witnessed significant changes in the
general conditions in which surveys have been undertaken, including
increased respondent resistance and sensitivity to response burden,
emphasis on voluntary surveys, increased incidence of proxy response,
relatively higher travel costs, and as a result increased and more
refined techniques of using telephone in conducting surveys, mail surveys,
etc. In such an environment, it is essential for a continuous survey to
maintain a capacity for testing various new options (or to study the
effects of changing conditions) with the view to developing and maintain-
ing a cost-efficient design. Basic requirements for such a capacity toge-
ther with the planning for a telephone experiment and its implications on
the current design are discussed in this section.

Lastly, some projects related to updating of the sample and develop-
ment of improved methodologies [ 16 ] are highlighted in the final section.

It should be noted that research into the areas of alternative sample
allocations and rotation patterns discussed in sections 2.2 and 2.3 are
in very initial stages. For that reason, the tables presented are the
results of preliminary investigations only, but they do indicate various
possibilities depending on the requirements and priorities of users for
the survey data.


## 2. DATA REQUIREMENTS

### 2.1 Establishment of Information Needs

It has now become a standard practice when designing a new survey or
redesigning an existing survey to determine the information needs
for establishing survey objectives at an early stage in the project.
For a continuous survey, failure to determine these needs on the assump-
tion that general objectives have remained unchanged would defeat one
of the most important purposes of the redesign exercise, namely eval-
uation of the survey from the viewpoint of its uses and effectiveness.

During the period between successive redesigns, information needs of users who participated in the earlier setting of objectives may have changed, and in addition new uses of the survey data other than those considered in the design of the survey may have emerged. It is incumbent upon the methodologists and the sponsors to discuss jointly with survey users the detailed information needs and priorities and to set up objecives for the redesign in clear and specific terms.

It is only when specific objectives have been agreed upon, that the design statisticians can properly discharge their responsibility of developing the most efficient survey design taking account of the operational constraints and the cost specified for the survey. The importance of this close dialogue at the initial stages cannot therefore be overemphasized. This is particularly so in the case of large-scale continuous surveys where major changes cannot usually be incorporated in midstream due to such factors as continuity of the time series, complexities of operations, and cost benefit considerations. Hence any failure at the outset in establishing survey objectives may continue to affect the survey results for the life of the design.

In an environment of fiscal restraint such as currently exists, a seemingly legitimate concern on the part of the survey sponsoring agency may be that initiation of discussions with the users might spark the type and degree of demands which could reach well beyond the scope of the survey. However,as long as the importance of budgetary constraints are clearly realized by both, there should be a definite advantage to such discussions. Not only would they enable the statisticians to take stock of and prioritize demands, but also they would serve to identify and inform users of those requirements which cannot be met by the survey and alert them to the consequences of misuse of the survey data. For fuller discussions on the role of user consultations in analysis of requirements, and on identification of feasibility, priority and method, we refer the reader to a paper by Fellegi and Ryten [ 5 ].

Another point to be emphasized is the lead time required in redesigning
a continuous survey and the importance of input from survey users
at the early stages of this process.  In contrast with adhoc
surveys where normally design activities may be completed in a
couple of months to about a year, the lead time needed for continuous
surveys is much longer.  By the same count, payoffs as well as stakes
are higher.  For the CLFS redesign, while the detailed research plans
are currently being formalized, some initial studies have already
been in progress since the middle of 1980 and the introduction of a
redesigned sample is scheduled for 1985.

As a beginning step in the process of determining information needs
for the coming CLFS redesign, members of a recently established
Evaluation Program For the LFS [ 15 ] will meet with all the major users
of CLFS data and a sample of other users, for the purposes of identifying
users' needs and how well the existing LFS satisfies these needs.
Based on these findings, the design statisticians will intensively
follow up cases where new information needs have emerged or earlier
information needs have changed. Success of the redesign program thus
becomes heavily dependent upon the timely specification of the require-
ments so as to provide the survey designers sufficient time to evaluate
alternatives and choose the most appropriate strategy for a given
situation.

The specification of information needs should include a specification
of;  characteristics of interest;  the types of estimates required-
rates, levels, changes in rates, or changes in levels; required fre-
quency of estimates; cross classifications (if any) desired for the
characteristic at different area levels of interest; and finally
associated data reliability requirements.

The specifications should also include a description of the uses of the
data, and their bearing on decision-making processes or allocation of
funds.  Equally essential is an assessment by the users of the importance

of the survey data for their program. Where information needs of a user are diverse, the user should also be asked to indicate priorities for them. Having received this input from users, overall priorities would be established and, subject to budgetary and other restrictions, would be translated into specific survey objectives. It should be emphasized that a primary responsibility of the design statistician in the process of user consultation is to provide assistance to the users in understanding what input is required of them and to provide technical guidance, for instance in the determination of reliability requirements, and identification of possible means of meeting their requirements.

It is worth drawing attention to a note by Platek [ 11 ], to papers by Cahoon, Kniceley and Shapiro [ 1 ], in which the importance of establishing clear survey ofjectives at an early stage of the survey has been emphasized.

To illustrate the importance of precise specification of survey objectives in deciding upon the choice of survey strategy, we present below three alternate means of meeting demands for more data. The choice of a particular strategy or combination of strategies should depend upon the type of data needed and the priorities set out for them. In the following sections alternatives will be presented with respect to reallocation of the sample (sect. 2.2), use of alternate rotation patterns (sect. 2.3) and lastly, expansion of the scope of the survey (sect. 2.4).


## 2.2   Sample Reallocation for Improved Monthly Subprovincial Estimates

Before discussing the implications of sample reallocation on data reliability, the expression for the coefficient of variation used in calculations is briefly discussed with relation to sample size, frequency of a characteristic, and design effect.

For the LFS, the coefficient of variation for monthly estimates for the characteristec unemployed ( u), and for an area ( a) of interest, can be expressed as

$$CV_a(u)\% = 100 \; \frac{SD_a(u)}{u_a}$$

where $\quad SD_a(u) = ( \displaystyle\sum_{t \varepsilon a} F_t \; (W_t - 1) \; P_t \; P_t \; q_t)$

where $\quad \displaystyle\sum_t$ = sum over strata or collection of strata for which sample design and sampling rate are the same.

and $\quad$ P = estimated persons 15 years of age or over

$\quad$ W = inverse sample rate (= $P/_n$, where n = sample size)

$\quad$ p = proportion unemployed

$\quad$ q = (1 - p)

$\quad$ F = design effect

$\quad U_a = \displaystyle\sum_{t \varepsilon a} P_t \; P_t$

From the above formulation, it can be seen that the reliability of estimates of level for a characteristic are primarily dependent on three factors:

i) Sample size: since W = $P/n$, other factors being constant, the CV% decreases proportionately to increases in the square root of the sample size. That is to reduce the CV in half, the sample size would have to be quadrupled.

ii) Frequency of the characteristic: the coefficient of variation is approximately inversely proportional to the square root of the proportion of persons having the characteristic. Thus for unemployed, the lower the unemployment rate, the larger the sample size required to obtain reliable estimates.

iii) Design effect: The design effect provides an overall comprehensive measure of the combined effect of all the design features, such as

stratification, multistage sampling and estimation. It is defined by the variance estimate obtained from the survey divided by the variance that would have resulted had the sample been taken in the form of a simple random sample of persons. The interpretation of a design effect of 2 for unemployed, would imply that, cost consideration aside, the sample design was only half as effective for measuring the characteristic unemployed as a simple random sample would have been. For the LFS, design effects are generally greater than one for most characteristics, due to the need for concentrating the sample in a relatively few selected areas as a means of reducing data collection costs.

Historically for the LFS, the characteristic unemployed has usually been considered of primary importance, and the total size and allocation of the sample have been determined to achieve specified reliability requirements for monthly estimates of unemployed. Prior to a sample size increase during the 1970's, the sample of 36,000 households was allocated with the primary objective of providing good monthly estimates for unemployed at the national level. When the sample was increased to 55,000 households, the additional sample was allocated on the basis of achieving more uniform reliability between provinces for unemployed. Because the increase was carried out after the redesign of the sample there was an additional restriction imposed by the sample design, namely that in Self Representing (SR) strata (i.e. larger cities) the sample could only be increased by multiples of the existing sample size, and in remaining (NSR) areas, increases had to be half-multiples of the existing sample size (i.e. 50%, 100%, 150%, etc.)

In increasing the sample size, uniform sampling rates by type of area (NSR and SR) within provinces were retained, as this provided an effective allocation scheme for improving provincial estimates for unemployed. Table 2.1 illustrates the impact of the increase on monthly CV's for the characteristic unemployed for the period Jan 75 to Dec 1980. The uniform sampling rates have the additional advantage of providing a good general purpose allocation considering the broad range of characteristics

on which information is collected, not only for the LFS, but by other surveys utilizing the LFS capacity.

Table 2.1      Pre-Increase and Post-Increase
CV% for Monthly Estimates of Unemployed

| Province | Post Increase Sample Size (households) (2) | % Increase (3) | CV% for unemployed | |
|---|---|---|---|---|
| | | | pre-increase (4) | post-increase (5) |
| Newfoundland | 3056 | 70.30 | 8.44 | 6.23 |
| Prince Edward Island | 1418 | 200.00 | 18.12 | 9.61 |
| Nova Scotia | 4021 | 29.80 | 6.55 | 5.34 |
| New Brunswick | 4217 | 67.78 | 8.23 | 5.44 |
| Quebec | 8541 | 17.06 | 4.56 | 3.54 |
| Ontario | 10850 | 14.24 | 4.31 | 3.65 |
| Manitoba | 4719 | 141.34 | 11.13 | 6.55 |
| Saskatchewan | 5724 | 200.00 | 14.42 | 7.56 |
| Alberta | 6709 | 100.00 | 8.65 | 6.49 |
| British Columbia | 6124 | 42.20 | 5.76 | 4.99 |
| Canada | 55379 | 52.12 | 2.33 | 1.88 |

With provincial CV's currently at acceptable levels, there has been an increased demand for more reliable subprovincial data. In the remainder of this section we examine how the reliability levels for subprovincial monthly estimates of unemployed could be improved by means of within province sample reallocations.

A disadvantage of the self-weighting design (uniform sampling ratio) is that for subprovincial regions variability in population sizes translate

into variations in the reliability of estimates. Currently monthly estimates of unemployed are published separately for 47 LFS Economic Regions for which the CV's are 25% or less. The remaining 19 ER's have been collapsed into groups of 2-4 to ensure that the reliability levels for the groups meet publication criteria.

It has been recently determined [ 17 ] that an additional sample of approximately 3000 dwellings would be required to achieve a 25% CV for each of the individual ER's where collapsing is carried out. It was also shown that these dwellings could be achieved by reallocating samples from the larger CMA's in the respective provinces. Refinements on such within province reallocations are currently being investigated using more months of survey data in the calculations and also taking into consideration reliabilities of estimates for other subprovincial areas as discussed in section 2.3 .

For illustration purposes, below we consider what could be achieved by within-province sample reallocations for the province of Manitoba. Present reliability levels for Manitoba's 8 Economic Regions based on data for the period Feb 78 to May 79 are shown in column 5 of Table 2.2. It might be noted that currently ER's 65 and 68, and ER's 61-64 are collapsed for publication purposes. The sample was reallocated to NSR and SR portions of individual Economic Regions so as to minimize data collection costs while achieving a fixed CV (22%) for unemployed following a general approach suggested by Fellegi et al [4], for all the ER's except 67 (Winnipeg), for which the sample size had to be reduced by 288 households. It should be noted that under the sample reallocation the provincial CV remains virtually unchanged, although costs would increase somewhat due to heavier sampling in NSR areas.

Table 2.2    Within Province Sample Reallocation for
Manitoba Economic Regions
(period Feb 78 - May 79)

| Economic Regions (1) | existing sample allocation | | | | reallocation of sample | | | |
|---|---|---|---|---|---|---|---|---|
| | $W_{NSR}$ (2) | $W_{SR}$ (3) | Hhlds (4) | CV(u)% (5) | $W'_{NSR}$ (6) | $W'_{SR}$ (7) | Hhlds (8) | CV(u)% (9) |
| 61 | 41.67 | - | 477 | 25.98 | 30.16 | - | 659 | 22.00 |
| 62 | 41.67 | - | 308 | 25.89 | 30.33 | - | 423 | 22.00 |
| 63 | 41.67 | 90.00 | 690 | 17.55 | 92.55 | 74.54 | 399 | 22.00 |
| 64 | 41.67 | 90.00 | 242 | 30.84 | 25.48 | 36.37 | 435 | 22.00 |
| 65 | 41.67 | - | 404 | 22.36 | 40.40 | - | 417 | 22.00 |
| 66 | 41.67 | 90.00 | 426 | 20.18 | 52.69 | 65.30 | 353 | 22.00 |
| 67 | - | 90.00 | 2030 | 7.73 | - | 104.88 | 1742 | 8.35 |
| 68 | 124.00[1] | 90.00 | 142 | 36.44 | 33.37 | 92.70 | 288 | 22.00 |
| Manitoba | 43.27 | 90.00 | 4719 | 6.15 | 39.37 | 100.52 | 4716 | 6.19 |

[1] remote area sample

There are some potential problems with an allocation scheme optimized
for the subprovincial estimates for the characteristec unemployed, that
have yet to be fully addressed, however.  For instance, it may be less
efficient for other surveys utilizing the LFS capacity.  While for other
surveys the desired allocations could be achieved by sub-sampling the LFS
selections, this would nevertheless reduce the sample size available to
such surveys.  Additionally the robustness of such an allocation against
changes in the unemployment levels over time would have to be studied
further.

## 2.3 Alternate Rotation Patterns

In a rotating panel survey such as the LFS, the monthly sample size
determines the reliability of monthly estimates of levels and rates;
however, it is primarily the rotation pattern which determines:
i) the reliability of estimates of change, whether month to month, quarter
   to quarter, or for a calendar month from one year to the next and
ii) the reliability of estimates obtained by combining monthly data to
   arrive at quarterly, semi-annual or annual estimates.
In general, rotation patterns which are better for i) are not as good
for ii) and vice versa. Thus the choice of a rotation pattern should
be governed by the relative priorities attached to these types of
estimates.

At the time of earlier redesigns of the LFS, there was little
demand for estimates of type ii) and therefore the choice of the
current LFS rotation pattern has reflected a predominant importance
for estimates of month to month change. Under the current LFS
rotation pattern, households remain in the sample for six consecutive
months, and each month one-sixth of the households rotate out of the
sample and are replaced by new ones. This scheme is very efficient
for measuring month to month changes as the 5/6 households in common
from one month to the next results in moderate to high correlations
between successive months' samples for most characteristics.

The same correlations between successive months' samples which are
advantageous for estimates of change are disadvantageous for average
estimates of level. As a result, the current LFS rotation pattern is
not as efficient for quarterly, semi-annual or annual estimates of rates
or level as some other schemes. It is of interest to compare the per-
formance of the LFS rotation pattern for combining data over months and
for estimates of month to month change with that of the Current Population
Survey (CPS), the counterpart of the LFS in the United States. In the

CPS households remain in the sample for 4 consecutive months, are out for 8 months, and then rotate back in again for 4 more months. Thus each month 1/4 of the households rotate.

If we denote $V_m$ as the variance of estimate for a given month m and $V_{cm}$ as the corresponding variance for estimates obtained by combining data for c months, then the variance reduction factor due to combining data, K, is defined as:

$$K = \frac{V_{cm}}{V_m}$$

Similarly, if we let $V_{(m, m+1)}$ denote the variance for estimates of change between months m and m+1, then the variance reduction factor for month to month change, $K'$, is defined as

$$K' = \frac{V_{(m, m+1)}}{V_m + V_{m+1}}$$

It should be noted that $K'$ is approximately equal to ( 1 - the correlation coefficient between the months' estimates). Table 2.3 presents values of K and $K'$ for the two rotation schemes for the characteristic unemployed. The smaller value for K for the CPS rotation scheme indicates that it is more efficient for combining data, while the smaller value of $K'$ for the LFS rotation scheme indicates it performs better for estimates of month to month change. It should be noted that the figures presented in the table for the LFS are the result of preliminary investigations only [ 7 ], and results for the CPS are taken from [ 18 ].

Table 2.3   Comparison of LFS and CPS Rotation Schemes for Unemployed

| | Variance reduction factor due to combining months data (K) | | | Variance reduction factor for month to month change (K') |
|---|---|---|---|---|
| | 3 mo | 6 mo | 12 mo | |
| CPS | .50 | .31 | .20 | .50 |
| LFS | .67 | .48 | .29 | .44 |

Research studies are planned to confirm the results of Table 2.3 for
the LFS, to consider similar variance factors for a broader range of
characteristics, and also to consider the implications on combined
estimates and estimates of change for other rotation patterns such as
3 - 9 - 3 (three months in the sample, 9 months out, and 3 months in) and
1 - 2 - 1 - 2 - 1, (one month in, 2 months out, one month in, etc.)

To further illustrate the impact of rotation pattern on average estimates
of level, Table 2.4 presents the sample sizes necessary to achieve 25%
CV's for annual estimates of unemployed for individual Census Divisions
for the LFS and CPS rotation patterns under two different allocation
schemes.

The augmentation allocation is based on retaining the present sample
allocation and adding to it whenever necessary to produce the required
reliability level for individual Census Divisions (CD's). The reallocat-
ion strategy on the other hand is based on a complete reallocation of
the sample to achieve the required reliability levels for CD's. On
practical considerations, both of these are extreme options and are used
only for illustrative purposes. The reallocation strategy in some cases
would result in a deterioration of monthly provincial estimates, while
the augmentation strategy is clearly too expensive. A comprehensive
strategy taking into consideration annual reliability levels for Census

Divisions, simultaneously with monthly reliability levels at the Province, Metropolitan Area, and Economic levels is currently under investigation.

Table 2.4          Additional Dwellings Required Monthly for
25% CV for Annual Estimates for
Unemployed for Census Divisions

| Province | LFS ROTATION PATTERN | | CPS ROTATION PATTERN | |
|---|---|---|---|---|
| | Augmentation | Reallocation | Augmentation | Reallocation |
| NFLS | -- | -- | -- | -- |
| PEI | -- | -- | -- | -- |
| NS | 90 | -- | 12 | -- |
| NB | 7 | -- | -- | -- |
| QUE | 2,520 | -- | 1,055 | -- |
| ONT | 2,300 | -- | 952 | -- |
| MAN | 3,578 | 1.411 | 1,887 | -- |
| SASK | 2,080 | -- | 877 | -- |
| ALTA | 1,357 | -- | 647 | -- |
| B.C. | 370 | -- | 184 | -- |
| Canada | 12,302 | 1,411 | 5,614 | -- |

It should be noted that the calculations are based on the assumption of the current LFS design and the unemployment level at the time of the 1976 Census. Further, assumed density factors of 1.5 and 3 are used for SR and NSR areas, and the variance reduction factors used are those given in table 2.3. Thus the figures in Table 2.4 should be considered only for the purpose of illustration and relative comparisons, as changes in any of the above factors including the design, will result in changed allocations.

The point clearly illustrated by Table 2.4 is that if sufficient priority
is attached to improved quarterly, semi-annual or annual average
estimates, then for cost-efficiency reasons, there would be a strong case
for changes in the rotation pattern. In this event, apart from the theo-
retical investigations of various rotation patterns, including the study of
the impact of rotation group biases on them, a detailed examination of
response burden and other operational aspects would have to be tested
through the field experimental capacity described in section 2.4. If,
on the other hand, higher priorities are given to the estimates of month
to month change, then the LFS rotation pattern should remain unchanged.


## 2.4  Current Survey Capacity

Recent years have witnessed an increased demand for more detailed labour
force data and data on a wide variety of characteristics influencing the
labour market situation. During the 1970's Statistics Canada successfully
responded to these demands by undertaking a major survey revision [ 10 ]
which included an expanded capacity for use of the LFS as a vehicle for
conducting other surveys. The current redesign will provide the oppor-
tunity to re-evaluate the role of the LFS in this regard.

Since the LFS is the only continuous household survey program carried
out by Statistics Canada, integration of other household surveys with
the LFS is desirable in the sense that these surveys can take advantage
of the investment the LFS represents in terms of sample frame,
design, data collection, and processing systems to obtain data more
quickly, at less cost and greater reliability than would be possible
through independent surveys. With the increased flexibility and capacity
of the LFS achieved through the revision and through methodological improve-
ments made at the last redesign, demands for use of the LFS as a vehicle
for conducting household surveys have continued to increase in recent years.
Examples of such surveys include: Survey of Consumer  Finances, Asset and
Debt, Family Expenditure, Annual Work Pattern, Household Facilities and

Equipment, Student Identification, Job Opportunity, Travel, Education, Smoking Habits, and Leisure Time Activities. Integration of these occasional surveys takes three different forms.

First, in the majority of cases, these surveys are conducted as supplements to the LFS due to cost and timeliness considerations. In such cases the most commonly adopted procedure is to collect data during the same visit, immediately after the LFS interview. In the case of enquiries with longer questionnaires, such methods of data collection as drop-off/pick-up are also utilized.

A second level at which the LFS has been utilized by other surveys is to select a different set of households in the same sampled areas as the LFS and to utilize labour force interviewers but at a different time period from the LFS. This is somewhat more costly than a supplement, but nevertheless represents a considerable saving over an independent survey. Examples of such use include; the Survey of Consumer Finances in odd years, when the survey content is expanded to include in depth questioning, for instance on Assets and Debts, and the program of Family Expenditure Surveys which consists of a recall survey and a diary survey.

The other situation in which LFS frame has been used is to select an 'independent sample' from the LFS frame, but in areas not currently being sampled by the LFS. The advantages over a totally independent sample are saving in sample design and implementation costs and also the control to avoid overlapping with the LFS and surveys associated with it. The Canada Health Survey for instance followed this approach in its survey design in cities, although in the remainder of the country a separate design was called for due to unique operational constraints.

Currently along with the LFS redesign activities, methodological aspects of other major surveys are also being researched. Just as it is important for the primary subject of enquiry using the continuing survey vehicle to

re-evaluate and re-establish its own objectives, it is equally incumbent on the other major users of the vehicle to follow the same course of action.

This will provide an opportunity for such surveys to maximize to the extent possible their benefits from the redesigned capacity of the vehicle, by determining optimal designs for their surveys, by being aware of implications of redesign alternatives, and by providing input to the decision processes.

Sponsors of each such major survey and the associated methodologists have recently begun this undertaking. Major studies in the optimization process would include stratification, sampling stages, allocation at various level of aggregations, determination of sampling and sub-sampling fractions, rotation patterns, response rates and their adjustments and other factors in the estimation process. It is not unlikely that these studies would result in a collection of optimal designs differing to a varying degree for different surveys.

Depending upon the importance attached to the major surveys using LFS vehicle and the degree to which the optimal designs differ, one of the three options may be followed, namely:
a) to redesign the current vehicle as a continuing household survey primarily for the LFS taking account of other surveys to the extent possible,
b) to redesign the current LFS vehicle as a general purpose survey or
c) to redesign the vehicle only for the LFS and develop a separate vehicle for conducting other major socio-economic surveys

The current situation is somewhere between a) and b); design features are optimized for the LFS, particularly in Non Self Representing Areas, (sect. 4.2); nevertheless the capacity is used extensively in a general purpose sense, as has already been described. It will be noticed that there is a very fine distinction between the options a) and b) and that

the difference mainly lies in the degree of importance associated to various subjects of enquiry using the redesigned vehicle.

In order to illustrate the distinction between the two approaches let us consider the problem of allocation of sample at a given level of aggregation (say R). Suppose there are m enquiries ($m = 1, 2, ..., M$) with the corresponding optimum allocation as $n_m$ at level R, with the LFS allocation being denoted by $n_{\ell}$. Say their magnitudes are as follows:

$$n_1 \leq n_2 \leq \cdots n_{\ell} \cdots \leq n_m \cdots \leq n_M$$

indicating that the subject M requires largest sample at level R. Note that this may happen at level R even if the total sample size for subject M may be smaller than that of the LFS.

In the case of option b) the approach would be to aim at a compromise allocation (say n*) such that $n_1 \leq n* \leq n_m$. In case of option a) however, the allocation would always be $n_{\ell}$ determined to be optimum for the LFS, and in order to accommodate other surveys, flexibility would have to be introduced into the vehicle to allow for over-sampling or sub-sampling as required. Option c) while having some technical merit, suffers from operational problems, such as co-ordination of the two vehicles to prevent overlapping samples. Even if such problems are taken care of, this option as such can be ruled out on the grounds of being very expensive unless some enquiries equally important and complex as the LFS come along.

Discussions and studies are being carried out in order to make a final decision on these options. It seems at this stage that the requirements of most surveys currently using the vehicle would be met under option a) by incorporating minor changes in various aspects of the LFS design and increasing the capacity of the survey vehicle as described below.

In order to meet the data requirements, studies are being undertaken
to develop an alternate small scale survey capacity in addition to
increasing the capacity for the current LFS vehicle. One component
of this program would be a collection of statistics on new subject
matters in anticipation of future requirements. This would thus
serve as a 'pilot' for full scale enquiries for more detailed data.
As well the small scale survey capacity would provide an opportunity
for analytical studies to examine inter-relationships between various
social and economic phenomena. The capacity may frequently be used
for surveys where it is necessary to react quickly in response to
data associated with policy concerns of the federal government. The
third area where this capacity would be useful is the development of
new techniques through well designed field experimentation. This
last aspect is discussed in more detail in the following section along
with data collection.

## 3. DATA COLLECTION

In a large scale survey, a single or a combination of data collection
methods such as personal interview, telephone and mail may be used,
depending upon the type of enquiry, available facilities, respondents'
attitudes, resource situation and timing constraints. Whatever be the method
adopted at the initial phase of a continuous survey, it requires regular
review as changes in the environment and conditions under which
data are collected will directly affect its quality. Over time, respon-
dents' attitudes towards surveys may change due to changing life style or
increased respondent burden; new tools and techniques may be developed;
legal requirements, the resource situation or quality of interviewers
may have changed. All these affect the quality of data collected and
hence the choice of method. Although certain changes in the data
collection procedure may be introduced at any time during the life of
the survey design, major changes affecting the cost and quality are
usually introduced along with the redesign of the vehicle. This is

because the procedure adopted for data collection affects both the
type of sampling design as well as the estimation procedure, and hence
to be cost effective the method of interview must be decided upon well
in advance of the sampling plan.  It should be noted that for a survey
vehicle like the LFS, which is used by various types of enquiries, the
effect of any change in the procedure of collecting data needs to be
investigated, including testing in the field, for as many of the major
enquiries as possible.  In this respect, just as in the case of establish-
ment of the survey objectives, close discussion and coordination among
sponsor, field staff and methodologist are important at the very early
stages of planning.

As mentioned in the previous section, a small scale capacity is being
considered to meet the current needs of social statistics.  One of the
primary purposes of establishing this capacity is to provide an oppor-
tunity for testing and developing new operational and methodological
procedures.  Testing of alternatives will focus on the data quality through
operational measures such as response rates, slippage, error rates, etc.,
and also the effect on the cost of the survey.  This methods test capacity
may be utilized depending on the purpose of the test in any of the follow-
ing manners:  use of same households as the LFS, separate set of households
in the same area as the LFS, or a completely different sample.  Also the
purpose of a particular test will determine its duration, location and
the spread of the sample.

It is expected that certain new methods and procedures will be tested
in the field with a view to examining suitability for the ongoing LFS.
As an example, one such test which deals with the extension of telephone
interviewing in the rural areas and smaller urban centres is briefly dis-
cussed below.

After a period of testing, the use of the telephone interview was expanded
at the time of the last redesign to cover all Self-Representing
Units primarily to reduce the cost of data collection. Currently
in the LFS, households are interviewed in person the first month they are
in the sample. In Self-Representing Units, if the respondent agrees to
the telephone, the interviews are as a rule conducted via telephone in
the second through sixth month the household remains in the sample. In
other areas interviews are conducted in person. A similar telephone
interviewing procedure will be tested for NSRU's. However, due to concern
over the confidentiality of the data, telephone interviewing will be restric-
ted to areas with a very low incidence of party lines.

Objectives of testing telephone interviewing in NSR areas will be to
determine for the LFS and other surveys using the vehicle:
i) Effect on data collection costs and sample design implications.
   Reductions in the travel component of collection costs and the
   potential for interviewers to handle larger assignments, could
   permit designs with less concentration of the sample, and hence a
   reduction in sampling variance. For instance, it might be possible
   to eliminate one or more stages of sampling.
ii) Data quality. Acceptance of telephone interviewing, effects on
   non response rates, and if possible on survey estimates would be
   examined.

The test would be conducted on a sub-set of the ongoing LFS inter-
viewer assignments, augmented in some cases by 10 - 20 percent to
study the effects of larger assignment sizes and different concen-
trations of the sample.

## 4. OTHER DESIGN RESEARCH

In this section, we briefly highlight some of the redesign projects
related to updating the sample and introduction of methodological
improvements in the sample design and estimation procedure.

## 4.1  Redesign of Self-Representing Units

Current LFS Self-Representing Units (SRU's) correspond to those cities
which were sufficiently large to yield a sample capable of supporting at
least one interviewer.  Minimum SRU sizes vary from a population of
10,000 in the Atlantic Region to 25,000 in Quebec and Ontario.  A first
step then will be a re-definition of the SR universe taking into consider-
ation impact of the 1976-77 sample size increase, population shifts, and
changes in boundaries of Census Metropolitan and Census Agglomeration
areas.

Larger SRU's are divided into sub-units and within sub-units, first
stage sampling units (i.e. clusters), are delineated on the basis of
field counts obtained in 1973.  The clusters correspond approximately
to city blocks.  A two stage sample of clusters, and dwellings (3 - 5
per selected cluster) is selected following a pps method based on
random groups of clusters [ 14 ] .  Using census data to simulate the
LFS design, research is being carried out to investigate the effects
on sampling effciency and operational suitability of alternative first
stage sampling units - such as census enumeration areas, blocks or
block faces - and of alternative allocations of the sample between and
within first stage units.

Another focus is on alternative means of achieving and maintaining an
up-to-date sample in SRU areas.  Due to the rapid and uneven growth which
occurs in these areas, without regular updating, the variance of estimates
can increase substantially [ 2 ].  Under the present sample updating
program [ 13 ], [ 3 ], for sub-units being updated, revised dwelling
counts for individual clusters are obtained on the basis of complete field
counts.  As an alternative to independently obtained field counts, the
use of census units, dwelling counts and maps in the redesign of the
sample is being investigated. Discussions are also in progress with Post

Canada concerning possible use of Post Canada maps and dwelling counts
to provide a future means of updating the LFS sample without incurring
the expense of field counts. The key to this would be the planned linking
of Postal Codes to 1981 Census units, and hence to LFS sampling units.

## 4.2 Redesign of Non Self-Representing Units

Non Self-Representing Units correspond to the smaller urban centres
and rural areas. In the present design, 1 - 5 geographically contiguous,
approximately equi-sized strata are formed within the NSR portions of
individual Economic Regions. Industry classifications were taken into
consideration in forming strata. Within strata, approximately 15 Primary
Sampling Units (PSU's) were delineated so as to be similar to the stratum
with respect to stratification variables and rural to urban population
ratios. To satisfy this last constraint, frequently urban centres had
to be shared amongst several PSU's within the stratum, often resulting
in discontiguity between rural and urban portions of PSU's.

Initially two PSU's per stratum were selected following the randomized
pps systematic method [ 8 ]. The sample was increased by selecting
additional PSU's [ 6 ], and at present 3 - 6 PSU's are selected per
stratum. It is felt that the sample increase strategy adopted may have
led to a reduction in the efficiency per unit cost of survey, although
the circumstances of the increase occuring in midstream ruled out more
technically desirable alternatives such as re-stratification to form an
increased number of strata, each with two selected PSU's.

As data requirements and design constraints, both technical and operational,
vary from province to province alternative designs will be investigated
by provinces or groups of provinces taken together as opposed to seeking
a uniform national design.

The NSR design is very much dependent not only on whether telephone interviewing is adopted as discussed in section 3, but also on the survey objectives. For instance, if an increased importance is attached to annual estimates for Census Divisions or to the estimates from other major surveys using the vehicle, then a design in which CD's were taken as primary strata would be seriously considered. In such a case, the design would likely feature rural/urban sub-stratification within CD's and utilization of either Census Sub Divisions or Census Enumeration Areas as first stage sampling units. Studies would be required to determine whether any loss in sampling efficiency for the LFS would be incurred under such a design, due to the reduced amount of optimal type stratification.

Additional studies in the NSR design which are planned, primarly to improve the design efficiency and facilitate updating include:

Buffer Areas

Generally growth in NSR areas is not large enough to warrant updating the sample between redesign. Exceptions, however, are the NSR areas close to the boundaries of certain Census Metreoplitan Areas. During the 10 year life cycle of the design, growth frequently reaches into these areas, where a more flexible design capable of being updated is therefore required.

Stages of Sampling

Studies will be conducted to determine the implications, both operational and theoretical of reducing the number of stages of sampling. This study would be closely linked to the study on telephone interviewing.

4.3  Estimation and Variance Estimation

A number of studies are planned into estimation and variance estimation procedures used by the LFS and other household surveys. Some of these

are briefly highlighted below:

## Final Ratio Adjustment

The current estimation procedure for individuals incorporates ratio estimation at the province level, using official population estimates by age-sex categories, adjusted for out of scope population (military and institutional). Research will be conducted into determining optimal age-sex post-strata, applying the ratio estimation at sub-provincial levels, and adjustment of LFS data for census undercount.

## Estimation for Family Units

In the past, post censal estimates of numbers of family units have been unavailable, with the result that there has been no standard procedure from one survey to the next for producing family based estimates. This project will address both of these problems, as well as attempting to ensure consistency between family and individual based estimates.

## Variance Estimation

Research will be carried out to compare alternative estimators with the Keyfitz [ 9 ] estimator currently being used from a point of view stability and extent of bias in the current estimator due to the violation of the sampling with replacement assumption. This will be examined for both seasonally adjusted as well as unadjusted sample estimates.

## Small Area Estimation

Research will continue into estimation methods for non-standard areas cutting across design strata. Estimators being studied include synthetic, composite, and sample regression. Attention is also being given to the treatment of large growth clusters falling into the sample, particularly as they affect estimates for small areas.

## 5. SUMMARY

While redesigning continuous surveys, the importance of close discussions among users, sponsor and designers at an early stage of the program is emphasized. This will not only help to re-evaluate the effectiveness of the ongoing program but it will be a useful exercise in identifying and informing the users about the limitations of the survey. As a result of such discussions the survey objectives can be established in the light of current and future data requirements. To illustrate the importance of the precise specification of the objectives, three alternate means of meeting data demands are discussed, namely reallocation, rotation patterns and survey capacity. Choice of these or other alternative would clearly depend on the specification of the information needs.

Like the specification of survey objectives, data collection procedure plays a very significant role in deciding the survey strategy for a particular situation. Designers aim at developing the most efficient design per unit cost and since the major part of the cost relates to the data collection, an early decision in this respect is essential. A small scale capacity is being developed to list and develop new procedures and it is planned to use this capacity in examining the suitability of telephone interviews in rural and smaller urban centres.

At present steps are being taken to establish the information needs and also to decide upon the field methodology. Several research and evaluation projects in the above context have been started. In addition research related to other aspects of the design and estimation methodology has begun.

RESUME

Les organismes spécialisés révisent périodiquement leurs enquêtes permanentes. Ces révisions tiennent à l'évolution des besoins en information auxquels l'enquête doit répondre et à l'evolution de la perception et de l'attitude du public à l'égard des enquêtes; elles sont analysées dans le contexte de la révision de l'enquête sur la population active du Canada après le recensement de 1981. En particulier, les auteurs font ressortir l'importance dès le début du processus de révision du dialogue entre les utilisateurs des données de l'enquête et les statisticiens concepteurs afin de déterminer les objectifs de l'enquête.

REFERENCES

[1] Cahoon, L.J., Kniceley, R.M. Jr. and Shapiro, G.M. (1980), "Informational Needs for Current Demographic Survey Design with Discussion of Key Redesign Research Projects", ASA Proceedings of the Section on Survey Research Methods, pp. 96-112.

[2] Drew, J.D., Choudhry, G.H. and Gray, G.B. (1978), "Some Methods for Updating Sample Survey Frames and their Effects on Estimation", ASA Proceedings of the Section on Survey Research Methods, pp. 62-71.

[3] Drew, J.D. and Singh, M.P. (1978), "An integrated Program for Continuous Redesign and Sample Size Reduction - A proposal for Self-Representing Areas", Technical Memorandum, Census and House-hold Survey Methods Division, Statistics Canada.

[4] Fellegi, I.P., Gray, G.B. and Platek, R. (1967), " The New Design of the Canadian Labour Force Survey", Journal of the American Statistical Association, 62, pp. 421-453.

[5] Fellegi, I.P., Ryten J. (1977), "An application of Functional Analysis - Current Trends in Statistics Canada", paper presented at the 25th Plenary Session of the Conference of European Statisticians.


[6] Gray, G.B. (1973), "On increasing the Sample Size (noⱼ of PSU's)", Internal Technical Memorandum, Census and Household Survey Methods Division, Statistics Canada.


[7] Gray, G.B. (1979), "Sampling Variance", Internal Technical Memorandum, Census and Household Survey Methods Division, Statistics Canada.


[8] Hartley, H.O. and Rao, J.N.K., (1962), "Sampling with Unequal Probabilities Without Replacement", Annals of Mathematical Statistics, 33, pp. 350-374.


[9] Keyfitz, N., (1957), "Estimates of Sampling Variance where Two Units are Selected from Each Stratum", Journal of the American Statistical Association, 52, pp. 503-510.


[10] Petrie, D.B., (1973), "Project Review: Assessment and Revision of the Canadian Labour Force Survey", discussion paper presented at Federal Provincial Conference.


[11] Platek, R. (1979), Lecture Notes for Survey Methodology Overview Course, Statistics Canada.

[12]   Platek, R. and Singh, M.P. (1976), "Methodology of the Canadian
       Labour Force Survey", Catalogue No. 71-526, Statistics Canada.


[13]   Platek, R. and Singh, M.P. (1978), "A Strategy for Updating
       Continuous Surveys", Metrika, Vol 25, pp. 1-7.


[14]   Rao, J.N.K., Hartley, H.O. and Cochran, W.F., (1962), "On a
       Simple Procedure of Unequal Probability Sampling Without Replace-
       ment", Journal of the Royal Statistical Society, 24(2), pp 482-490.


[15]   Sayant, G., Pold, H. and Macredie, I., (1981), "Project Proposal -
       Program Evaluation Labour Force Survey - Summary", Internal
       Report, Statistics Canada.


[16]   Singh, M.P., and Drew, J.D., (1981), Research Plans for the
       Redesign of the Canadian Labour Force Survey", for presentation
       at 1981 ASA Meetings, Section on Survey Research Methods.


[17]   Tarte, F., (1981), "Improvement of Reliability for Economic
       Regions", Internal Technical Memorandum, Census and Household
       Survey Methods Division, Statistics Canada.


[18]   U.S. Department of Commerce, Bureau of the Census (1978), "The
       Current Population Survey Design and Methodology", Technical Paper
       40, Washington, D.C., U.S. Government Printing Office, pp. 64, 96-99.

FOR-HIRE TRUCKING SURVEY:

SURVEY DESIGN

R. Lussier[1]

The methodology of the For-hire Trucking Survey is discussed
in this paper.  This survey provides good examples of
administrative and operational constraints faced by survey
statisticians and field data collection teams.

## 1.  INTRODUCTION

This paper presents the methodology of the For-hire Trucking Survey, a
multi-stage probability survey of shipping documents retained by for-hire
trucking carriers in Canada.  The paper is structured as follows.  In the
second section, the survey context is described.  In the third section,
the ultimate sampling unit, namely the shipment is defined.  Then the
universe and the frame of the survey are depicted in the fourth section.
After that, some major administrative considerations are listed in the
fifth section and the stratification and sample allocation are covered
in the sixth section.  Then the first stage sample design and the
subsequent stage(s) sample design are presented respectively in the
seventh and eighth sections.  In the ninth section, the field operations
are discussed.  The data processing and the estimation methods
respectively are explained in the tenth and eleventh sections.  Finally,
a comment about the future of the survey concludes the paper in the
twelfth section.

---

[1]  R. Lussier, Business Survey Methods Division, Statistics Canada.

## 2. SURVEY CONTEXT

### 2.1 Primary Objective and Uses of the Survey

The primary objective of this survey is to provide information about
the domestic intercity movements of goods by the Canadian for-hire
trucking industry. The fo-hire trucking industry covers any carrier
which for compensation undertakes the transport of goods by truck. This
This survey measures the output of this industry in terms of revenues
earned, tons carried and ton-miles performed by commodity group.

Requests for estimates from this survey come from a wide variety of
sources such as government departments concerned with trade; transport
regulatory officials at both federal and provincial levels; carriers;
university consultants; industry associations; and many other organi-
zations and individuals who share a common interest in transportation.

The estimates are used extensively to serve four basic requirements.
First, they measure the volume of domestic trade transported by inter-
city for-hire carriers provincially and interprovincially. Secondly,
they provide a cross-check on the rate of industrial growth reflected
by intercity commodity movements and they provide information on
regional development. Thirdly, they assist in transportation studies
(e.g. [1], [2] and [3]). Finally, they support the presentation of
briefs, submissions and other inquiries to regulatory authorities and
commissions.

### 2.2 Background

Initial work on the For-hire Trucking Survey began in 1969. At that
time, a study of various methods of collecting commodity origin and
destination statistics was carried out. The study results showed that,
from a cost-benefit point of view, a sample survey of shipping documents
was the only viable approach to collect the data.

In 1970, a pilot survey was conducted to assess the survey approach effectiveness. The pilot survey involved the examination of the shipping documents of about 200 for-hire trucking firms throughout the country. The favourable response to the pilot survey and the availability of origin, destination, commodity, weight and revenue information of the shipping records indicated that the survey approach was feasible.

The For-hire Trucking Survey has been conducted on an annual basis since 1970 by the Transportation and Communications Division of Statistics Canada. However, the survey design has changed over time. Examples of changes mentioned in this paper are changes to the frame and changes to the sample allocation technique.


## 3. THE ULTIMATE SAMPLING UNIT

The 1969 study and the 1970 pilot survey mentioned in Section 2 indicated that the ultimate sampling unit be the shipment.

The principal characteristics needed from each sampled shipment are the true origin and the final destination; the description of the commodity(ies) carried; the weight; the transportation revenue earned and the interlined shipment information. An interlined shipment occurs when a consignment is moved by a carrier to an intermediate point and then moved by another carrier to another point. The interlined shipment information is used to unduplicate interlined shipments.

The secondary characteristics needed are the month and year of shipment; the quantity of commodity and the unit of measurement (e.g. 5 board feet, 20 gallons, 15 sacks, etc.); the method of movement (e.g. heated van, refrigerated van, piggyback, fishyback, container, etc.); some remarks on shipment weight transcribed (e.g. minimum weight, convenient weight used for calculating revenue, etc.); the rate charged and the rate

condition codes (e.g. code indicating where rate is minimum, per 100 lb., per hour, etc.); and the revenue condition codes (e.g. code indicating where exact transportation revenue is not available, where shipment is out of scope, etc.).

The shipment characteristics can be found on documents known in the trucking industry as Probills, Bills of Lading, Load Manifests, Trip Reports, Invoices, or a combination of the above in either a computer listing format or other media of storage.

These documents can be filed in complete numeric sequence; in broken numeric sequence; in chronological order; in alphabetical order (e.g. by customer name); by terminal; by commodity type, i.e. usually contracts; or in no order at all. The documents may even be cross-filed; for example, by serial number and by customer's name. Within a filing system, documents may be kept in a set of file drawers, in a set of binders or shannon files, on shelves, in drawers, or even in a book.

## 4. UNIVERSE AND FRAME

The first choice for the frame is ideally a list of all shipments. However, such a list is not available. Instead, D.S.L.P.'s (Document Storage Location Points) are used as natural clusters of shipments for the first stage sampling units of the design. A D.S.L.P. is a site at which shipping documents suitable for sampling are kept. A trucking company may have more than one D.S.L.P. in the case where shipping documents are stored at several terminals but not at the company head office. In some cases, a D.S.L.P. contains shipping documents for more than one company.

The universe and the frame of D.S.L.P.'s have changed since the survey started in 1970.

In 1970, the universe was defined to include the D.S.L.P.'s of all provincially regulated or licensed carriers, regardless of size, type or major activity. The frame for this universe was derived from provincial license lists and included about 15,000 D.S.L.P's.

Since 1975, the universe has been limited to the D.S.L.P's of the carriers earning $100,000 or more annually from intercity trucking. Other exclusions are mentioned in the survey publication [4]. Also, the frame for a given year has been the list of the D.S.L.P.'s of the carriers whose reports to the Motor Carriers Freight and Household Goods Movers Survey for the previous year show earnings of $100,000 or more from the domestic movement of goods over more than 15 miles of public roads.

The description of the frame of the 1978 survey and of the principal statistics estimated from the survey data are given in Table 1 at the end of the paper.

## 5. ADMINISTRATIVE RESTRICTIONS

The survey design consists of selecting shipments from the files of selected D.S.L.P.'s and of transcribing the characteristics of the selected shipments on coding sheets. However, administrative restrictions limit the number of transcriptions. The adminstrative restrictions are presented in this section and the survey design is detailed in the next sections.

### 5.1  Maximum Total Number of Transcriptions

The cost of the survey is relatively high. As an example, the 1977 For-hire Trucking Survey cost $494,000. The distribution of this expenditure by function is given in Table 2 at the end of the paper.

The cost has set the maximum total number of transcriptions to 225,000 shipments. This size has remained the same since 1972 although the

number of shipments carried by the for-hire trucking industry has
increased.

## 5.2  Minimum Number of Transcriptions per Selected D.S.L.P.

A minimum number of shipments are to be transcribed from the filing
systems of the D.S.L.P.'s in the sample to justify travel and salary
expenditures.  Under the present design, a minimum of 200 shipments are
selected from each D.S.L.P. in the sample.

## 5.3  Maximum Number of Transcriptions per Selected D.S.L.P.

Identification of shipment records in the sample and the transcription
of information from these records are done at the D.S.L.P.'s.  There is
a constraint on the number of days the data collection team spends at
a particular location, so that the respondents are not burdened by the
presence of the team.  This constraint translates to a maximum of 3000
shipment records to be transcribed from any one D.S.L.P.

## 6.  STRATIFICATION AND SAMPLE ALLOCATION

Using the results of the previous year's Motor Carriers Freight and House-
hold Goods Movers Survey, the D.S.L.P.'s are stratified according to their
intercity transportation revenue class, their type of operation and their
area of operation.  The intercity transportation revenue class indicates if
the D.S.L.P. earned $2 million or more, between $500,000 and $1,999,999,
or between $100,000 and $499,999 dollars of revenue from intercity
freight transport.  The type of operation says if the D.S.L.P. is a
general freight carrier, an automobile carrier, a household goods mover,
a van line, a bulk (e.g. petroleum, milk, etc.) carrier or an other
specialized (e.g. heavy machinery, livestock, etc.) carrier.  The area
of operation is different depending on the D.S.L.P. total transportation
revenue.  If the revenue is greater than or equal to 2 million, the area
of operation indicates if the revenue is between 2 and 20 million or

greater than 20 million dollars and if the predominant source of revenue is earned east of the Manitoba/Ontario border, west or from international trucking. If the revenue is less than 2 million dollars, the area of operation indicates which of the 10 provinces, the Yukon, the Northwest Territories or the international brought the predominant source of revenue to the D.S.L.P. There were 102 non-empty strata in the 1978 For-hire Trucking Survey.

Once the frame of D.S.L.P.'s is stratified, the number of shipments to be selected and transcribed ultimately from each stratum is determined by allocating the maximum total number of transcriptions to strata. The allocation technique has changed since the survey started in 1970.

Originally, the allocation was calculated in three steps. First, the total number of transcriptions was allocated to 5 domains so that the coefficients of variation were equal for ton-miles over the 5 domains. These domains were the geographic regions of origin of the shipments. The coefficients of variation were estimated using historical data. This first step gave a number of transcriptions to each domain. Secondly, the number of transcriptions to each domain was allocated to strata using essentially a Neyman allocation. This second step gave a number of transcriptions to each domain within each stratum. Finally, the numbers of transcriptions to each domain within each stratum were summed over the domains to get the total stratum allocation of transcriptions.

In 1975, the allocation scheme was revised because it was felt that the estimates of the true variance of each domain within each stratum needed for the Neyman allocation were not reliable and because the resultant strata allocations were highly variable across years and adversely affected longitudinal analyses.

The revised allocation procedure is radically different than the original one. It has been developed using years of experience and is partially a

judgement allocation based on results of the previous year's survey. It consists of allocating workloads defined as 100 transcriptions. The allocation is performed in three steps.

First, the total number of transcriptions, i.e. 2250 workloads, is allocated to the groups of strata having the same intercity transportation revenue class as follows:

| Intercity Transportation Revenue Class | Workload Allocation to Groups of Strata |
|---|---|
| $2 million or more | 908 |
| $500,000 to $1,999,999 | 832 |
| $100,000 to $499,999 | 510 |
| TOTAL | 2,250 |

Secondly, the allocation of workloads to strata within a group of strata having the same intercity transportation revenue class is performed using a stratum size measure. This measure for a stratum is the total intercity transportation revenue in units of 10,000 dollars of the D.S.L.P.'s in the stratum. For the group of strata having intercity transportation revenue of $2 million or more, the allocation is proportional to stratum size measure. For the others, the allocation is proportional to the square root of the stratum size measure. The square root is used in the latter case because otherwise some strata having little contribution to the revenue would almost be ignored in the sample.

Finally, the allocations obtained from the previous steps are adjusted as follows. The allocation is reduced to 2 workloads in the strata of the international carriers and household goods movers with intercity transportation revenue of $2 million or more. It is increased in strata

where detailed data are needed. It is also adjusted to meet the administrative restrictions mentioned in Section 5 and to preserve consistency with previous years' allocations.


## 7. FIRST STAGE SAMPLE DESIGN


The current first stage consists of selecting a twice replicated stratified sample of D.S.L.P.'s. The sample selection is different in class 1 strata than in the other strata. Class 1 strata cover class 1 D.S.L.P.'s which are D.S.L.P.'s earning 2 million or more dollars of intercity transportation revenue. The two sample selection procedures are described below. The selection is done by methodologists of Business Survey Methods Division.


### 7.1 Selection of D.S.L.P.'s in Class 1 Strata and Allocation of Workloads to Selected D.S.L.P.'s.


All class 1 D.S.L.P.'s are selected with probability one. The reason for this approach is that these D.S.L.P.'s are know to be heterogeneous with respect to the principal statistics estimated.

Each class 1 D.S.L.P. next must be assigned a number of workloads for each replicate of the sample. This assignment has to be derived from the stratum allocation which was obtained through the procedure described in Section 6. The distribution of the stratum allocation to individual D.S.L.P.'s is done as follows. Let the stratum allocation be x workloads and let the number of D.S.L.P.'s in the stratum be d. One workload is assigned per replicate for each D.S.L.P. in the stratum so a total of 2d workloads are assigned. The remaining $w = (x - 2d)$ workloads are equally distributed to the two replicates. Then a probability proportional to size systematic sample of $(\frac{w}{2})$ D.S.L.P.'s is

drawn for each replicate. The measure of size used in the intercity transportation revenue. One workload is assigned per selection to the selected D.S.L.P.

## 7.2 Selection of D.S.L.P.'s in the Other Strata and the Allocation of Workloads to Selected D.S.L.P.'s

The selection of D.S.L.P.'s in other strata and the allocation of workloads to selected D.S.L.P.'s is done simultaneously as follows. Let the stratum allocation be k workloads. Each replicate in the stratum is assigned $(\frac{k}{4})$ workloads. A probability proportional to size systematic sample of $(\frac{k}{4})$ D.S.L.P.'s is drawn for each replicate. The measure of size used is the intercity transportation revenue. One workload is assigned per selection to the selected D.S.L.P. Finally, the workload assignments to selected D.S.L.P.'s are doubled.

## 7.3 Review of the Sample

The selected D.S.L.P.'s and the distribution of workloads are reviewed using the information from the previous survey, and adjustments to assignments are made, if warranted for practical reasons.

# 8. SUBSEQUENT STAGE(S) SAMPLE DESIGN

The subsequent stage(s) of the sample design consist(s) of selecting shipments from the files of each selected D.S.L.P. This selection is done by Statistics Canada Regional Operations Division staff at the D.S.L.P. The sample design is different depending on whether the filing system of the D.S.L.P. is small or large.

If the filing system is small (i.e. less than or equal to 2,000 shipments), then two independent systematic samples of 50 shipments are selected for each workload assigned to the D.S.L.P..

If the filing system is large, (i.e. greater than 2,000 shipments), then a sample of about 100 shipments is selected independently for each workload in two stages. For the first stage, an estimate of the number of bundles (defined as 100 shipments) is obtained and divided into 8 equal sections. Then a bundle is selected at random from each section and the selected bundle is located in the filing system. For the second stage, a systematic sample of shipments is selected from the selected bundle. The sample interval used for the systematic sampling is 8 so that 12 or 13 shipments are usually transcribed. Thus, the 8 sections provide approximately 100 transcriptions for each workload.

## 9.  FIELD OPERATIONS

This section discusses the activities that involve the Statistics Canada Regional Operations staff: namely the training of the Regional Operations project managers, the planning of the collection, the collection itself and some special cases.

### 9.1  Training of the Regional Operation Project Managers

Every year, the Statistics Canada Regional Operations project managers are trained on all aspects of the survey. The training is five days long and is conducted during the month of March. The survey project manager as well as methodologists are involved in the training. A collection procedures manual is used for the in-class training. The Regional Operations project managers also receive on the job training by visiting a number of D.S.L.P.'s with different filing systems.

### 9.2  Planning of the Collection

Every spring, the Regional Operations project managers recruit the interviewers and administer a thorough training program. Then the interviewers with the advice of their Regional Operations project manager schedule their work, plan their itineraries and telephone D.S.L.P

officials for appointments. The collection takes place between May and September for the survey covering the shipments of the previous calendar year.

## 9.3 Overall Description of the Collection

When the interviewer gets on the D.S.L.P. premises, he/she has first to conduct an interview with the D.S.L.P. officials. During the interview, he/she will explain the survey, will mention the uses of the data, will estimate the time required to do the work and will complete a control form. The control form records information about the operations of the firm such as the total tonnage transported; the total number of shipments carried; the types of commodities carried and the percentage each type represents in the total transportation revenue; and the filing system(s) used.

Often, the interviewer has a choice of filing systems which provide information on the items needed in this survey. The interviewer assesses the completeness of information on principal characteristics from various filing systems, and then chooses the most appropriate system.

Then the interviewer estimates the number of shipments in the selected filing system of the carriers. This estimation is needed to be able to properly select a sample of shipments and to calculate the weights of the sampled shipments. This estimation involves measurement when the filing system is neither numeric nor broken numeric.

After then, the interviewer selects the sample shipments and transcribes their characteristics. This latter operation is often difficult because it can be hard to understand the various documents and the coding used on some documents. The interviewer often has to interpret the information on the documents and to enter on the coding sheets the data in a format that would be accepted by the computer system.

Finally, the survey project manager and the methodologists are consulted whenever necessary to assist in the field operations.

## 9.4 Special Cases

This sub-section discusses synthesising, abortions and cancellations.

Synthesising is the construction of hypothetical workloads when a D.S.L.P. does not keep documents suitable for sampling. In such a case, the interviewer collects through an interview with the D.S.L.P. officials macro-information about the D.S.L.P. Then this information is sent to Transportation and Communications Division who constructs shipment data in a format accepted by the computer system.

Abortions are in-scope D.S.L.P.'s for which we have not obtained transcriptions nor macro-information. Examples of abortions are a D.S.L.P which refuses to cooperate or a D.S.L.P. on strike. The contributions of the abortions are reflected in the estimates via imputations using previous year's data or via adjustments to the weights.

Finally, cancellations are D.S.L.P.'s which are identified as out-of-scope for the reference period. In spite of the efforts made in verifying the D.S.L.P.'s in the universe against several sources on the activity of the carriers, interviewers find out that some D.S.L.P.'s in the sample are out-of-scope for the reference period. In such cases, Head Office makes adjustments to the weights when a large number of cancellations occur within a stratum.

## 10. DATA PROCESSING

All data processing is carried out at Head Office. Incoming data first go through a manual edit procedure which uses the data collected on the control forms. If the shipment data are correct and complete, they are sent to Key-edit. A standard data record file is created by matching the incoming transcriptions and accompanying material with a check-in list called workload master file. Out-of-scope shipments are discarded. There

were about 59,000 out-of-scope shipments in the 1978 For-hire Trucking Survey, i.e. about 26.2% of the transcribed shipments were out-of-scope. Some types of out-of-scope shipments are shipments to or from the U.S.A.; shipments transported 15 miles or less from origin to destination; shipments which were off-highway; shipments which would be double counted as a result of interlining between road carriers; shipments which would be double counted because they were recorded by household goods movers who are van line agents and by the van lines themselves; shipments which relate to a period other than the reference period; shipments which did not bring any intercity transportation revenue; and records which relate to non-transportation services such as storage, packing, equipment rental, labour loading and unloading.

The in-scope shipment records are assigned Standard Geographic Codes, Standard Commodity Codes and the distance between the origin and destination of the shipment using respective computer libraries.

Missing fields and those failing edits are imputed. There are two major imputation procedures that may take place in the system. These are prorating which is arithmetic imputation using fixed relationships and simple Hot Deck which consists of matching the record with missing data to a "similar" record with complete data. Essentially these procedures take values (or codes) to be imputed from valid or complete records and are applied to records which are incomplete. These procedures are premised on the assumption that the characteristics of records within the same workload are similar.

Weights are finally assiged and the data file is passed to the estimation module.

Detailed diagnostics produced at each stage of data processing are used as a quality check on the data passing through the system.

## 11. ESTIMATION METHOD

The estimates are generated from a very small sample relative to the
size of the population. As an example, the 1978 sample represented only
about 0.5% of the shipments in-scope to the survey.

Each workload generates an independent estimate for its stratum as follows.
First, using the estimated size of the filing system in shipments, the
individual shipment data of the workload are expanded to obtain an estimate
at the D.S.L.P. level. Next, the relative size of the D.S.L.P. to the
stratum size in terms of intercity transportation revenue is used to expand
the D.S.L.P. estimate to the stratum level.

The average of estimates from all workloads in a given replicate within
a given stratum provides the replicate stratum estimate. These replicate
estimates are averaged to derive an overall stratum estimate. Finally,
these overall stratum estimates are aggregated to provide national
estimates.

Standard errors of the estimates are calculated using the two replicated
estimates for each stratum.

The tabulated estimates are reviewed to ensure a check on the accuracy
of the weights as well as a general check on the quality of the estimates
generated. Examples of the review are the comparisons of the estimates
to the previous year estimates from the For-hire Trucking Survey and to
the previous year estimates from the Motor Carrier Freight Survey.

## 12. FUTURE OF THE SURVEY

The survey is being re-designed. One of the objectives is to develop a
new approach to the data collection to reduce the cost of collection per
in-scope sampled shipment so that the sample size can be enlarged. The

possibility of using computer tapes as a vehicle for the reporting firms falls into this approach. The progress of the re-design project is being monitored by an interdepartmental committee involving Transport Canada, the Canadian Transport Commission and Statistics Canada.

## ACKNOWLEDGEMENT

## RESUME

Cet article décrit la méthodologie de l'enquête sur le Transport routier de marchandises pour le compte d'autrui. Cette enquête fournit de bons exemples de contraintes administratives et opérationnelles rencontrées par les statisticiens d'enquêtes et par les équipes de collecte de données sur le terrain.

## REFERENCES

[1] Cairns, M. and Kirk, B.D. (1980), "Canadian For-hire Trucking and the Effects of Regulation: A cost Structure Analysis", Canadian Transport Commission, Research Report No. 10-80-03.

[2] Chow, G. (1980), "An Analysis of Selected Aspects of Performances of For-hire Motor Carriers in Canada", report prepared for the Bureau of Competition Policy, Consumer and Corporate Affairs Canada.

[3] McRae, J.R. and Prescott, D.M. (1980), "Definition and Characteristics of the Trucking Markets: A Statistical Analysis", report prepared for Transport Canada.

[4]   Statistics Canada, "For-hire Trucking Survey", Catalogue 53-224,
      Annual.

TABLE  1

Description of the Frame of the 1978 Survey and of the Principal
Statistics Estimated from the Survey Data by Intercity Revenue Group.

| Statistics | Intercity Revenue | | | |
| --- | --- | --- | --- | --- |
| | $2,000,000 and more | $500,000 to $1,999,999 | $100,000 to $499,999 | Total |
| Number of D.S.L.P.'s in the frame | 218 | 540 | 1,339 | 2,097 |
| Sample allocation in shipments | 90,800 | 83,200 | 51,000 | 225,000 |
| Number of shipments used in tabulation | 69,620 | 62,729 | 33,646 | 165,995 |
| Estimated number of shipments in-scope to this survey | 20,146,157 | 7,016,909 | 5,261,362 | 32,424,428 |
| Estimated revenue ($,000) | 1,720,578 | 449,781 | 303,575 | 2,473,934 |
| Estimated number of tons  (,000) | 62,703 | 36,743 | 22,800 | 122,426 |
| Estimated number of ton-miles (,000) | 16,594,065 | 5,277,815 | 3,120,894 | 24,992,774 |

TABLE  2

1977 For-hire Trucking Survey Expenditure by Functions

| Function | Expenditure |
|---|---|
| Field Data Collection | $295,000 |
| Mangement, Operations | $ 98,000 |
| Data Processing | $ 65,000 |
| Methodological Support | $ 30,000 |
| Travel, Printing, Misc. | 6,000 |
| Total | $494,000 |

# CONSTRUCTION OF WORKING PROBABILITIES AND
# JOINT SELECTION PROBABILITIES FOR FELLEGI'S
# PPS SAMPLING SCHEME

### G.H. Choudhry[1]

A FORTRAN Subroutine to obtain the "working probabilities" for Fellegi's (1963) method of unequal probability sampling is given. The solution is obtained by an iterative procedure where the starting values for the (k+1)th draw "working probabilities" are the solutions for the kth draw "working probabilities" and the iterative procedure is terminated when a prespecified accuracy is achieved. The limitation is that the Subroutine can only be used to obtain upto and including the 5th draw "working probabilities". It was observed that the convergence occurs very fast in double precision. Therefore all real variables have been declared as double precision. The joint selection probabilities $\Pi_{ij}$'s i.e. the probability that both the ith and jth units are in the sample are obtained by summing the probabilities of selecting those samples that contain both the ith and jth units. The joint selection probabilities are required for the variance estimation of the Horvitz-Thompson estimator of population total of the characteristic of interest.

## 1. DESCRIPTION

Fellegi (1963) has proposed a method for selecting a sample of n ($\geq 2$) units draw by draw and without replacement out of N units in such a way that the probability for the i-th unit to be selected is equal to $p_i$ at each of the n successive draws ( $\sum_{i=1}^{N} p_i = 1$ ). This is achieved by determining (n-1) sets of selection probabilities referred to as "working probabilities". Let the (n-1) sets of "working probabilities" be

---

[1]  G.H. Choudhry, Census and Household Survey Methods Division, Statistics Canada.

$$p_i(k) \geq 0, \quad i = 1, 2, \ldots, N; \quad k = 2, 3, \ldots, n$$

$$\sum_{i=1}^{N} p_i(k) = 1, \quad k = 2, 3, \ldots, n.$$

The $p_i(k)$, $i = 2, 3, \ldots, N$ are the "working probabilities" for selecting a unit at the k-th draw. The selection probabilities at the first draw $p_i(1)$ are given by

$$p_i(1) = p_i, \quad i = 1, 2, \ldots, N.$$

Then the overall (unconditional) probability $\delta_i(k)$ of selecting i-th unit at the k-th draw is given by

$$\delta_i(k) = \sum_{(k-1;\ i)} \left[ p_{i_1}(1) \times \frac{p_{i_2}(2)}{1 - p_{i_1}(2)} \times \ldots \times \frac{p_{i_{k-1}}(k-1)}{1 - p_{i_1}(k-1) - p_{i_2}(k-1) \ldots - p_{i_{k-2}}(k-1)} \right.$$

$$\left. \times \frac{p_i(k)}{1 - p_{i_1}(k) - p_{i_2}(k) \ldots - p_{i_{k-1}}(k)} \right]$$

$$i = 1, 2, \ldots, N;$$

$$k = 1, 2, \ldots, n$$

where $\sum_{(k-1;\ i)}$ denotes the summation over all possible ordered $(k-1)$ - tuples of $(i_1, i_2, \ldots, i_{k-1})$ such that $i_1, i_2, \ldots, i_{k-1}$ are different integers between 1 and N, and none of them is equal to i. The condition that the i-th unit be selected with probability $p_i$ at each of the n successive draws is satisfied by setting

$$\delta_i(k) = p_i, \quad i = 1, 2, \ldots, N; \quad k = 1, 2, \ldots, n.$$

We have $p_i(1) = p_i$, $i=1, 2, \ldots, N$. Given that $p_i(2), \ldots, p_i(k-1)$ have already been found, then approximate $p_i^{(0)}(k)$ by $p_i(k-1)$ and obtain $p_i^{(1)}(k)$ from the following formula

$$p_i^{(m)}(k) = p_i \times \left\{ \sum_{(k-1;\ i)} \left[ p_{i_1}(1) \times \frac{p_{i_2}(2)}{1-p_{i_1}(2)} \times \ldots \times \frac{p_{i_{k-1}}(k-1)}{1-p_{i_1}(k-1)-p_{i_2}(k-1)\ldots-p_{i_k}(k-1)} \right. \right.$$

$$\left. \left. \times \frac{1}{1-p_{i_1}^{(m-1)}(k)-p_{i_2}^{(m-1)}(k) \ldots - p_{i_k}^{(m-1)}(k)} \right] \right\}^{-1}$$

by setting $m = 1$ for $i = 1, 2, \ldots N$. Repeat for $m = 2, 3, \ldots$, etc. until $p_i^{(m)}(k) = p_i^{(m-1)}(k)$ for all $i$ up to the required number of decimal places. The procedure is carried out for $k = 2, 3, \ldots n$, thus obtaining the $(n-1)$ sets of "working probabilities" $p_i(2)$, $p_i(3), \ldots, p_i(n)$. Since $i$-th unit is selected with probability equal to $p_i$ at each of the $n$ successive draw, this property of the scheme makes it very attractive for rotating sample designs.

Bayless and Rao (1977) excluded Felleig's (1963) method from their study for n=4 due to covergence problems with the routine they used for obtaining the "working probabilities". They were not getting satisfactory answers even after a large number of iterations especially when c.v.(x)[*] was not small, where x-values are the sizes of the units in the population.

---

[*] $c.v.(x) = \dfrac{(\sum\limits_{i=1}^{N} x_i^2 - (\sum\limits_{i=1}^{N} x_i)^2 /N)/ (N-1)}{(\sum\limits_{i=1}^{N} x_i /N)}$

We have used Fellegi's (1963) example for which C.V. (x) is small and two
populations [Cochran (1978) and Kish (1965)] with larger values for
C.V.(x) to obtain the "working probabilities" for selecting upto 4 units.
The iterative procedure was terminated when the change between two
successive iterations was less than $10^{-6}$ for each element of the solution
vector. The description of the populations and the number of iterations
require to obtain the "working probabilities" at each of the draws is
given below:

| Pop. No. | Source | N | C.V.(X) | No. of iterations at draw | | |
|------|--------|---|---------|---|---|---|
| | | | | 2 | 3 | 4 |
| 1 | Fellegi [1963, p. 198] | 6 | 0.25 | 5 | 7 | 12 |
| 2 | Cochran [1978, p. 152] | 20 | 1.03 | 4 | 5 | 7 |
| 3 | Kish [1965, p. 42 ] | 20 | 1.19 | 4 | 6 | 8 |

It is noticed that for the three populations we have used, the convergence
at each of the draws is obtained in a very few number of iterations
although the number of iterations required at each successive draw
increases. It should be remarked that the values of "working probabi-
lities" obtained for Fellegi's (1963) example agree with his values.

The joint selection probabilities are required for estimating the
variance of the Horvitz-Thompson estimator

$$\hat{Y} = \frac{1}{n} \sum_{i \varepsilon s} \frac{y_i}{p_i}$$

of the total $Y = \sum_{i=1}^{N} y_i$ of y - variate of interest, where $y_i$ is the value

of y - variate pertaining to the i-th unit. Let $\pi_{ij}$ denote the probability

that both the i-th and j-th units are included in the sample, then $\pi_{ij}$,

i=1, 2, ... N-1; j= i + 1, i + 2, ..., N can be obtained as follows:

Let $\delta_{ij}$ (k, $\ell$) denote the probability that the i-th unit was selected at

the k-th draw and the j-th unit was selected at the $\ell$th draw ($\ell > k$). The

probability $\delta_{ij}$ (k, $\ell$) is given by:

$$\delta_{ij}(k,\ell) = \sum_{(\ell-2;i,j)} [p_{i_1}(1) \times \frac{p_{i_2}(2)}{1-p_{i_1}(2)} \times \dots \times \frac{p_{i_{k-1}}(k-1)}{1-p_{i_1}(k-1)-p_{i_2}(k-1)\dots -p_{i_{k-1}}(k-1)} \times$$

$$\frac{p_i(k)}{1-p_{i_1}(k)-p_{i_2}(k) \dots - p_{i_{k-1}}(k)} \times \frac{p_{i_{k+1}}(k+1)}{1-p_{i_1}(k+1)-p_{i_2}(k+1) \dots -p_{i_{k-1}}(k+1)-p_i(k+1)} \times$$

$$\dots \times \frac{p_{i_{\ell-1}}(\ell-1)}{1-p_{i_1}(\ell-1)-p_{i_2}(\ell-1) \dots -p_{i_{k-1}}(\ell-1)-p_i(\ell-1)-p_{i_{k+1}}(\ell-1) \dots p_{i_{\ell-2}}(\ell-1)} \times$$

$$\frac{p_j(\ell)}{1-p_{i_1}(\ell) - p_{i_2}(\ell) \dots - p_{i_{k-1}}(\ell) - p_i(\ell) - p_{i_{k+1}}(\ell) \dots - p_{i_{\ell-1}}(\ell)} ],$$

$$i \neq j = 1, 2, \dots, N;$$

$$k = 1, 2, \dots, n-1;$$

$$\ell = k+1, k+2, \dots, n$$

where $\sum\limits_{(\ell-2;i,j)}$ denotes the summation over all possible ordered $(\ell-2)$-tuples

of $(i_1, i_2,\ldots,i_{k-1}, i_{k+1}, \ldots i_{\ell-1})$ such that $i_1, i_2,\ldots,i_{k-1}, i_{k+1},\ldots,i_{\ell-1}$

are different integers between 1 and N, and none of them is equal to i or j.

Then $\pi_{ij}$, the probability that both i-th and j-th units are included in the

sample, is given by

$$\pi_{ij} = \sum_{k=1}^{n-1} \sum_{\ell=k+1}^{n} [\delta_{ij}(k,\ell) + \delta_{ji}(k,\ell)],$$

$$i = 1, 2, \ldots, N-1;$$

$$j = i+1, i+2, \ldots, N.$$

## Structure

SUBROUTINE  WKPROB (N, NS, MA, P, P1, P2, P3, P4, Q1, Q2, DEL,

MAX, ACC, PI, TOL, IFAULT)

Formal parameters - all real parameters in double precision.

| N | Integer | Input: | number of units in the population |
|---|---------|--------|-----------------------------------|
| NS | Integer | Input: | sample size, $2 \le NS \le 5$ |
| MA | Integer | Input: | dimension of PI in the calling program |
| P | Real Array(N) | Input: | contains the relative measure of sizes of units in the sequence $P_1, P_2, \ldots, P_N$; |

$$\sum_{i=1}^{N} P_i = 1$$

| P1 | Real Array (N) | Output: | working probabilities for selecting a unit at the 2nd draw |
|----|----------------|---------|----------------------------------------------------------|
| P2 | Real Array (N) | Output: | working probabilities for selecting a unit at the 3rd draw |
| P3 | Real Array (N) | Output: | working probabilities for selecting a unit at the 4th draw |

| | | | |
|---|---|---|---|
| P4 | Real Array(N) | Output: | working probabilities for selecting a unit at the 5th draw. |
| Q1 | Real Array(N) | Workspace | |
| Q2 | Real Array(N) | Workspace | |
| DEL | Real Array (MA,NS) | Workspace | |
| MAX | Integer | Input: | maximum number of interations allowed for obtaining each set of working probabilities |
| ACS | Real | Input: | desired accuracy of the working probabilities |
| P1 | Real Array (MA,MA) | Output: | matrix returning the joint selection probabilities $\pi_{ij}$, $i = 1, 2, \ldots, N-1$ $j = i + 1, i + 2, \ldots N$ |
| TOL | Real | Input: | maximum allowed value for the absolute difference between $\sum_{i=1}^{N} p_i$ and the number one |
| IFAULT | Integer | Output: | failure indicator |

Failure Indications

IFAULT = 0   normal termination

= 1   one or more of $p_i > (1/NS)$

= 2   $DABS(\sum_{i=1}^{N} p_i - 1.0) > TOL$

= 3   both conditions 1 and 2 occur

= 4   sample size greater than 5

= 5   desired accuracy was not obtained in maximum allowed number of interations

## ACKNOWLEDGEMENT

RESUME

L'auteur expose un sous-programme FORTRAN visant à obtenir
les "probabilités de travail" à l'aide de la méthode
d'échantillonnage à probabilités inégales de Fellegi (1963).
On obtient la solution par une méthode itérative dans
laquelle les valeurs de départ des "probabilités de travail"
du (k   1)-ième tirage sont la solution du k-ième tirage des
"probabilités de travail"; ce calcul prend fin lorsque l'on
atteint un niveau de précision déterminé à l'avance.  Le
sous-programme est limité car son utilisation ne peut
dépasser le 5$^e$ tirage des "probabilités de travail".  On a
observé que la convergence se produit très rapidement en
double précision.  Par conséquent, toutes les variables
réelles ont été déclarées en double précision.  Les prob-
abilités conjointes de sélection, c.-à-d. la probabilité que
les i-ième et j-ième unités fassent toutes deux partie de
l'échantillon, s'obtiennent par sommation des probabilités
de sélection des échantillons contenant les deux unités en
cause.  Les probabilités conjointes de sélection sont
nécessaires à l'estimation de la variance de l'estimateur
Horvitz-Thompson du total de la caractéristique à l'étude
dans la population.

REFERENCES

[1]  Bayless, D.L. and Rao, J.N.K. (1970), "An empirical study of
     Stabilities of Estimators and Variance Estimators in Unequal
     Probability Sampling (n=3 or 4)", Journal of the American
     Statistical Association, 65, 1645-1667.

[2]  Cochran, W.G. (1977),Sampling Techniques, 3rd Ed., New York:
     John Wiley and Sons.

[3]  Fellegi, I.P. (1963), "Sampling with and without Replacement:
     Rotating and Non-Rotating Samples", Journal of the American
     Statistical Association, 58, 183-201.

[4]  Kish, L. (1965), Survey Sampling, New York:  John Wiley and Sons.

```
C.....SUBROUTINE TO OBTAIN WORKING PROBABILITIES AND            00001000
C.....JOINT SELECTION PROBABILITIES FOR FELLEGI'S PPS           00002000
C.....SAMPLING SCHEME. REF: 1963 JASA 58 , PP 183-201 .         00003000
C.....                                                          00004000
      SUBROUTINE WKPROB(N,NS,MA,P,P1,P2,P3,P4,Q1,Q2,DEL,MAX,ACC, 00005000
     1             PI,TOL,IFAULT)                                00006000
      IMPLICIT REAL*8(A-H,O-Z)                                  00007000
      DIMENSION P(N),P1(N),P2(N),P3(N),P4(N),Q1(N),Q2(N),DEL(MA,NS), 00008000
     1             PI(MA,MA)                                     00009000
C.....N     IS POPULATION SIZE.                                 00010000
C.....NS    IS SAMPLE SIZE AND CAN HAVE VALUES 2 , 3 , 4 , 5 .  00011000
C.....MA    IS MAXIMUM DIMENSION  OF  PI IN THE MAIN PROGRAM.   00012000
C.....P     IS THE VECTOR OF  GIVEN PROBABILITIES . SUM P = 1.0 00013000
C.....P , P1 , P2 , P3 , P4 AND P5 ARE SELECTION PROBABILITIES AT 00014000
C.....1ST , 2ND , 3RD , 4-TH AND 5-TH DRAWS RESPECTIVELY .      00015000
C.....Q1 , Q2 , DEL .... WORK SPACE .                           00016000
      CONTINUE                                                  00017000
C.....MAX IS THE MAXIMUM NUMBER OF ITERATIONS ALLOWED TO OBTAIN 00018000
C.....THE WORKING PROBABILITIES.                                00019000
C.....ACC IS THE DESIRED ACCURACY OF THE WORKING PROBABILITIES. 00020000
C.....PI IS THE OUTPUT RETURNING THE JOINT SELECTION PROBABLITIES . 00021000
C.....TOL IS THE PARAMETER SO THAT SUM P CANNOT DEVIATE FROM 1.0 BY 00022000
C.....MORE THAN THE VALUE OF TOL.                               00023000
C.....IFAULT IS  FAILURE INDICATOR TAKING THE FOLLOWING VALUES: 00024000
C.....                                                          00025000
C..... 0            IF PI COMPUTED, NORMAL TERMINATION.         00026000
C..... 1            IF NS*P .GE.1.0   FOR ONE OR MORE P VALUES . 00027000
C..... 2            IF DABS(SUM P - 1.0) IS GREATER THAN TOL.   00028000
C..... 3            IF BOTH OF THE ABOVE TWO CONDITIONS.        00029000
C..... 4            IF NS , THE SAMPLE SIZE, IS GREATER THAN 5 . 00030000
C..... 5            IF DESIRED ACCURACY NOT OBTAINED IN MAXIMUM 00031000
C.....              ALLOWED NUMBER OF ITERATIONS.               00032000
C.....                                                          00033000
      IFAULT=4                                                  00034000
      IF(NS.GT.5) RETURN                                        00035000
      IFAUL1=0                                                  00036000
      IFAUL2=0                                                  00037000
      IDRAW=1                                                   00038000
      XNS=NS                                                    00039000
      SUMP=0.0                                                  00040000
      DO 1 I=1,N                                                00041000
      SUMP=SUMP+P(I)                                            00042000
      Q1(I)=P(I)                                                00043000
      DEL(I,1)=P(I)                                             00044000
      IF(XNS*P(I).GT.1.0) IFAUL1=1                              00045000
    1 CONTINUE                                                  00046000
      IF(DABS(SUMP-1.0).GT.TOL) IFAUL2=2                        00047000
      IFAULT=IFAUL1+IFAUL2                                      00048000
      IF(IFAULT.NE.0) RETURN                                    00049000
C.....                                                          00050000
```

```
C.....SELECTING UNIT 2 .                                             00051000
C.....                                                               00052000
      IDRAW=IDRAW+1                                                  00053000
      A=0.0                                                          00054000
      DO 20 J=1,N                                                    00055000
   20 A=A+F0(N,J,P,Q1)                                              00056000
      ICOUNT=0                                                       00057000
   21 ICOUNT=ICOUNT+1                                                00058000
      IF(ICOUNT.GT.MAX) GO TO 999                                   00059000
      DMAX=0.0                                                       00060000
      DO 22 I=1,N                                                    00061000
      DEN=A-F0(N,I,P,Q1)                                            00062000
      Q2(I)=P(I)/DEN                                                 00063000
      DIFF=DABS(Q2(I)-Q1(I))                                        00064000
      IF(DIFF.GT.DMAX) DMAX=DIFF                                    00065000
      Q1(I)=Q2(I)                                                   00066000
      A=DEN+F0(N,I,P,Q1)                                            00067000
   22 CONTINUE                                                       00068000
      IF(DMAX.GT.ACC) GO TO 21                                      00069000
      WRITE(3,24) IDRAW,ICOUNT                                      00070000
   24 FORMAT(1H1,////,20X,'WORKING PROBABILITIES AT DRAW       : ',I5, 00071000
     1          ////,20X,'NUMBER OF ITERATIONS FOR CONVERGENCE = ',I6, 00072000
     2          ////)                                               00073000
      DO 25 I=1,N                                                    00074000
      P1(I)=Q1(I)                                                   00075000
      DEL(I,2)=P1(I)                                                 00076000
      WRITE(3,26) I,P1(I)                                           00077000
   25 CONTINUE                                                       00078000
   26 FORMAT(1H0,20X,' PROB ( ',I2,' ) = ',D14.6)                  00079000
      IF(IDRAW.EQ.NS) GO TO 550                                     00080000
C.....                                                               00081000
C.....SELECTING UNIT 3.                                              00082000
C.....                                                               00083000
      IDRAW=IDRAW+1                                                  00084000
      A=0.0                                                          00085000
      DO 30 J=1,N                                                    00086000
      DO 30 K=1,N                                                    00087000
   30 A=A+F1(N,J,K,P,P1,Q1)                                         00088000
      ICOUNT=0                                                       00089000
   37 ICOUNT=ICOUNT+1                                                00090000
      IF(ICOUNT.GT.MAX) GO TO 999                                   00091000
      DMAX=0.0                                                       00092000
      DO 31 I=1,N                                                    00093000
      S1=0.0                                                         00094000
      DO 32 J=1,N                                                    00095000
      S1=S1+F1(N,I,J,P,P1,Q1)+F1(N,J,J,P,P1,Q1)+F1(N,J,I,P,P1,Q1)   00096000
   32 CONTINUE                                                       00097000
      DEN=A-S1+2.0*F1(N,I,I,P,P1,Q1)                                00098000
      Q2(I)=P(I)/DEN                                                 00099000
      DIFF=DABS(Q2(I)-Q1(I))                                        00100000
```

```
      IF(DIFF.GT.DMAX) DMAX=DIFF                                   00101000
      Q1(I)=Q2(I)                                                  00102000
      S1=0.0                                                       00103000
      DO 33 J=1,N                                                  00104000
      S1=S1+F1(N,I,J,P,P1,Q1)+F1(N,J,J,P,P1,Q1)+F1(N,J,I,P,P1,Q1)  00105000
   33 CONTINUE                                                     00106000
      A=DEN+S1-2.0*F1(N,I,I,P,P1,Q1)                               00107000
   31 CONTINUE                                                     00108000
      IF(DMAX.GT.ACC) GO TO 37                                     00109000
      WRITE(3,24) IDRAW,ICOUNT                                     00110000
      DO 36 I=1,N                                                  00111000
      P2(I)=Q1(I)                                                  00112000
      DEL(I,3)=P2(I)                                               00113000
      WRITE(3,26) I,P2(I)                                          00114000
   36 CONTINUE                                                     00115000
      IF(IDRAW.EQ.NS) GO TO 550                                    00116000
C.....                                                            00117000
C.....SELECTING UNIT 4.                                           00118000
C.....                                                            00119000
      IDRAW=IDRAW+1                                                00120000
      A=0.0                                                        00121000
      DO 40 J=1,N                                                  00122000
      DO 40 K=1,N                                                  00123000
      DO 40 L=1,N                                                  00124000
   40 A=A+F2(N,J,K,L,P,P1,P2,Q1)                                   00125000
      ICOUNT=0                                                     00126000
   49 ICOUNT=ICOUNT+1                                              00127000
      IF(ICOUNT.GT.MAX) GO TO 999                                  00128000
      DMAX=0.0                                                     00129000
      DO 41 I=1,N                                                  00130000
      S1=0.0                                                       00131000
      S2=0.0                                                       00132000
      DO 42 J=1,N                                                  00133000
      DO 43 K=1,N                                                  00134000
      S1=S1+F2(N,I,J,K,P,P1,P2,Q1)+F2(N,J,J,K,P,P1,P2,Q1)          00135000
     1      +F2(N,J,I,K,P,P1,P2,Q1)+F2(N,J,K,I,P,P1,P2,Q1)         00136000
     2      +F2(N,J,K,J,P,P1,P2,Q1)+F2(N,J,K,K,P,P1,P2,Q1)         00137000
   43 CONTINUE                                                     00138000
      S2=S2+2.0*F2(N,I,I,J,P,P1,P2,Q1)+F2(N,J,J,I,P,P1,P2,Q1)      00139000
     1      +2.0*F2(N,I,J,I,P,P1,P2,Q1)+F2(N,J,I,J,P,P1,P2,Q1)     00140000
     2      +2.0*F2(N,J,I,I,P,P1,P2,Q1)+F2(N,I,J,J,P,P1,P2,Q1)     00141000
     3      +2.0*F2(N,J,J,J,P,P1,P2,Q1)                            00142000
   42 CONTINUE                                                     00143000
      DEN=A-S1+S2-6.0*F2(N,I,I,I,P,P1,P2,Q1)                       00144000
      Q2(I)=P(I)/DEN                                               00145000
      DIFF=DABS(Q2(I)-Q1(I))                                       00146000
      IF(DIFF.GT.DMAX) DMAX=DIFF                                   00147000
      Q1(I)=Q2(I)                                                  00148000
      S1=0.0                                                       00149000
      S2=0.0                                                       00150000
```

```
      DO 44 J=1,N                                              00151000
      DO 45 K=1,N                                              00152000
      S1=S1+F2(N,I,J,K,P,P1,P2,Q1)+F2(N,J,J,K,P,P1,P2,Q1)      00153000
     1    +F2(N,J,I,K,P,P1,P2,Q1)+F2(N,J,K,I,P,P1,P2,Q1)       00154000
     2    +F2(N,J,K,J,P,P1,P2,Q1)+F2(N,J,K,K,P,P1,P2,Q1)       00155000
   45 CONTINUE                                                 00156000
      S2=S2+2.0*F2(N,I,I,J,P,P1,P2,Q1)+F2(N,J,J,I,P,P1,P2,Q1)  00157000
     1    +2.0*F2(N,I,J,I,P,P1,P2,Q1)+F2(N,J,I,J,P,P1,P2,Q1)   00158000
     2    +2.0*F2(N,J,I,I,P,P1,P2,Q1)+F2(N,I,J,J,P,P1,P2,Q1)   00159000
     3    +2.0*F2(N,J,J,J,P,P1,P2,Q1)                          00160000
   44 CONTINUE                                                 00161000
      A=DEN+S1-S2+6.0*F2(N,I,I,I,P,P1,P2,Q1)                   00162000
   41 CONTINUE                                                 00163000
      IF(DMAX.GT.ACC) GO TO 49                                 00164000
      WRITE(3,24) IDRAW,ICOUNT                                 00165000
      DO 47 I=1,N                                              00166000
      P3(I)=Q1(I)                                              00167000
      DEL(I,4)=P3(I)                                           00168000
      WRITE(3,26) I,P3(I)                                      00169000
   47 CONTINUE                                                 00170000
      IF(IDRAW.EQ.NS) GO TO 550                                00171000
C.....                                                         00172000
C.....SELECTING UNIT 5.                                        00173000
C.....                                                         00174000
      IDRAW=IDRAW+1                                            00175000
      A=0.0                                                    00176000
      DO 50 J=1,N                                              00177000
      DO 50 K=1,N                                              00178000
      DO 50 L=1,N                                              00179000
      DO 50 M=1,N                                              00180000
   50 A=A+F3(N,J,K,L,M,P,P1,P2,P3,Q1)                          00181000
      ICOUNT=0                                                 00182000
   59 ICOUNT=ICOUNT+1                                          00183000
      IF(ICOUNT.GT.MAX) GO TO 999                              00184000
      DMAX=0.0                                                 00185000
      DO 51 I=1,N                                              00186000
      S1=0.0                                                   00187000
      S2=0.0                                                   00188000
      S3=0.0                                                   00189000
      DO 52 J=1,N                                              00190000
      DO 53 K=1,N                                              00191000
      DO 54 L=1,N                                              00192000
      S1=S1+F3(N,J,K,I,L,P,P1,P2,P3,Q1)+F3(N,J,K,J,L,P,P1,P2,P3,Q1)  00193000
     1    +F3(N,J,K,K,L,P,P1,P2,P3,Q1)+F3(N,I,J,K,L,P,P1,P2,P3,Q1)   00194000
     2    +F3(N,J,J,K,L,P,P1,P2,P3,Q1)+F3(N,J,I,K,L,P,P1,P2,P3,Q1)   00195000
     3    +F3(N,J,K,L,I,P,P1,P2,P3,Q1)+F3(N,J,K,L,J,P,P1,P2,P3,Q1)   00196000
     4    +F3(N,J,K,L,K,P,P1,P2,P3,Q1)+F3(N,J,K,L,L,P,P1,P2,P3,Q1)   00197000
   54 CONTINUE                                                 00198000
      S2=S2+2.0*F3(N,I,J,I,K,P,P1,P2,P3,Q1)+F3(N,J,J,I,K,P,P1,P2,P3,Q1) 00199000
     1    +2.0*F3(N,J,I,I,K,P,P1,P2,P3,Q1)+F3(N,J,I,J,K,P,P1,P2,P3,Q1)  00200000
```

```
      2        +2.0*F3(N,J,J,J,K,P,P1,P2,P3,Q1)+F3(N,I,J,J,K,P,P1,P2,P3,Q1) 00201000
      3        +2.0*F3(N,I,I,J,K,P,P1,P2,P3,Q1)+F3(N,J,K,J,I,P,P1,P2,P3,Q1) 00202000
      4        +2.0*F3(N,J,K,I,I,P,P1,P2,P3,Q1)+F3(N,J,K,K,I,P,P1,P2,P3,Q1) 00203000
      5        +2.0*F3(N,I,J,K,I,P,P1,P2,P3,Q1)+F3(N,J,J,K,I,P,P1,P2,P3,Q1) 00204000
      6        +2.0*F3(N,J,I,K,I,P,P1,P2,P3,Q1)+F3(N,J,K,I,J,P,P1,P2,P3,Q1) 00205000
      7        +2.0*F3(N,J,K,J,J,P,P1,P2,P3,Q1)+F3(N,J,K,K,J,P,P1,P2,P3,Q1) 00206000
      8        +2.0*F3(N,J,J,K,J,P,P1,P2,P3,Q1)+F3(N,J,I,K,J,P,P1,P2,P3,Q1) 00207000
      9        +    F3(N,J,K,I,K,P,P1,P2,P3,Q1)+F3(N,J,K,J,K,P,P1,P2,P3,Q1) 00208000
      A        +2.0*F3(N,J,K,K,K,P,P1,P2,P3,Q1)+F3(N,I,J,K,J,P,P1,P2,P3,Q1) 00209000
      B        +    F3(N,I,J,K,K,P,P1,P2,P3,Q1)+F3(N,J,J,K,K,P,P1,P2,P3,Q1) 00210000
      C        +    F3(N,J,I,K,K,P,P1,P2,P3,Q1)                             00211000
   53 CONTINUE                                                             00212000
      S3=S3+6.0*F3(N,I,I,I,J,P,P1,P2,P3,Q1)                                00213000
      1        +6.0*F3(N,I,J,I,I,P,P1,P2,P3,Q1)                            00214000
      2        +2.0*F3(N,J,J,I,I,P,P1,P2,P3,Q1)                            00215000
      3        +6.0*F3(N,J,I,I,I,P,P1,P2,P3,Q1)                            00216000
      4        +2.0*F3(N,J,J,J,I,P,P1,P2,P3,Q1)                            00217000
      5        +2.0*F3(N,J,I,J,I,P,P1,P2,P3,Q1)                            00218000
      6        +2.0*F3(N,I,J,J,I,P,P1,P2,P3,Q1)                            00219000
      7        +6.0*F3(N,I,I,J,I,P,P1,P2,P3,Q1)                            00220000
      8        +2.0*F3(N,J,J,I,J,P,P1,P2,P3,Q1)                            00221000
      9        +2.0*F3(N,J,I,I,J,P,P1,P2,P3,Q1)                            00222000
      A        +6.0*F3(N,J,J,J,J,P,P1,P2,P3,Q1)                            00223000
      B        +2.0*F3(N,J,I,J,J,P,P1,P2,P3,Q1)                            00224000
      C        +2.0*F3(N,I,J,I,J,P,P1,P2,P3,Q1)                            00225000
      D        +2.0*F3(N,I,J,J,J,P,P1,P2,P3,Q1)                            00226000
      E        +2.0*F3(N,I,I,J,J,P,P1,P2,P3,Q1)                            00227000
   52 CONTINUE                                                             00228000
      DEN=A-S1+S2-S3+24.0*F3(N,I,I,I,I,P,P1,P2,P3,Q1)                      00229000
      Q2(I)=P(I)/DEN                                                       00230000
      DIFF=DABS(Q2(I)-Q1(I))                                              00231000
      IF(DIFF.GT.DMAX) DMAX=DIFF                                          00232000
      Q1(I)=Q2(I)                                                          00233000
      S1=0.0                                                               00234000
      S2=0.0                                                               00235000
      S3=0.0                                                               00236000
      DO 55 J=1,N                                                          00237000
      DO 56 K=1,N                                                          00238000
      DO 57 L=1,N                                                          00239000
      S1=S1+F3(N,J,K,I,L,P,P1,P2,P3,Q1)+F3(N,J,K,J,L,P,P1,P2,P3,Q1)       00240000
      1        +F3(N,J,K,K,L,P,P1,P2,P3,Q1)+F3(N,I,J,K,L,P,P1,P2,P3,Q1)   00241000
      2        +F3(N,J,J,K,L,P,P1,P2,P3,Q1)+F3(N,J,I,K,L,P,P1,P2,P3,Q1)   00242000
      3        +F3(N,J,K,L,I,P,P1,P2,P3,Q1)+F3(N,J,K,L,J,P,P1,P2,P3,Q1)   00243000
      4        +F3(N,J,K,L,K,P,P1,P2,P3,Q1)+F3(N,J,K,L,L,P,P1,P2,P3,Q1)   00244000
   57 CONTINUE                                                             00245000
      S2=S2+2.0*F3(N,I,J,I,K,P,P1,P2,P3,Q1)+F3(N,J,J,I,K,P,P1,P2,P3,Q1)   00246000
      1        +2.0*F3(N,J,I,I,K,P,P1,P2,P3,Q1)+F3(N,J,I,J,K,P,P1,P2,P3,Q1) 00247000
      2        +2.0*F3(N,J,J,J,K,P,P1,P2,P3,Q1)+F3(N,I,J,J,K,P,P1,P2,P3,Q1) 00248000
      3        +2.0*F3(N,I,I,J,K,P,P1,P2,P3,Q1)+F3(N,J,K,J,I,P,P1,P2,P3,Q1) 00249000
      4        +2.0*F3(N,J,K,I,I,P,P1,P2,P3,Q1)+F3(N,J,K,K,I,P,P1,P2,P3,Q1) 00250000
```
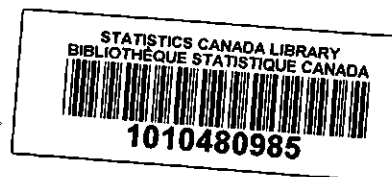
```
      5      +2.0*F3(N,I,J,K,I,P,P1,P2,P3,Q1)+F3(N,J,J,K,I,P,P1,P2,P3,Q1)  00251000
      6      +2.0*F3(N,J,I,K,I,P,P1,P2,P3,Q1)+F3(N,J,K,I,J,P,P1,P2,P3,Q1)  00252000
      7      +2.0*F3(N,J,K,J,J,P,P1,P2,P3,Q1)+F3(N,J,K,K,J,P,P1,P2,P3,Q1)  00253000
      8      +2.0*F3(N,J,J,K,J,P,P1,P2,P3,Q1)+F3(N,J,I,K,J,P,P1,P2,P3,Q1)  00254000
      9      +    F3(N,J,K,I,K,P,P1,P2,P3,Q1)+F3(N,J,K,J,K,P,P1,P2,P3,Q1)  00255000
      A      +2.0*F3(N,J,K,K,K,P,P1,P2,P3,Q1)+F3(N,I,J,K,J,P,P1,P2,P3,Q1)  00256000
      B      +    F3(N,I,J,K,K,P,P1,P2,P3,Q1)+F3(N,J,J,K,K,P,P1,P2,P3,Q1)  00257000
      C      +    F3(N,J,I,K,K,P,P1,P2,P3,Q1)                              00258000
   56 CONTINUE                                                             00259000
      S3=S3+6.0*F3(N,I,I,I,J,P,P1,P2,P3,Q1)                                00260000
      1      +6.0*F3(N,I,J,I,I,P,P1,P2,P3,Q1)                              00261000
      2      +2.0*F3(N,J,J,I,I,P,P1,P2,P3,Q1)                              00262000
      3      +6.0*F3(N,J,I,I,I,P,P1,P2,P3,Q1)                              00263000
      4      +2.0*F3(N,J,J,J,I,P,P1,P2,P3,Q1)                              00264000
      5      +2.0*F3(N,J,I,J,I,P,P1,P2,P3,Q1)                              00265000
      6      +2.0*F3(N,I,J,J,I,P,P1,P2,P3,Q1)                              00266000
      7      +6.0*F3(N,I,I,J,I,P,P1,P2,P3,Q1)                              00267000
      8      +2.0*F3(N,J,J,I,J,P,P1,P2,P3,Q1)                              00268000
      9      +2.0*F3(N,J,I,I,J,P,P1,P2,P3,Q1)                              00269000
      A      +6.0*F3(N,J,J,J,J,P,P1,P2,P3,Q1)                              00270000
      B      +2.0*F3(N,J,I,J,J,P,P1,P2,P3,Q1)                              00271000
      C      +2.0*F3(N,I,J,I,J,P,P1,P2,P3,Q1)                              00272000
      D      +2.0*F3(N,I,J,J,J,P,P1,P2,P3,Q1)                              00273000
      E      +2.0*F3(N,I,I,J,J,P,P1,P2,P3,Q1)                              00274000
   55 CONTINUE                                                             00275000
      A=DEN+S1-S2+S3-24.0*F3(N,I,I,I,I,P,P1,P2,P3,Q1)                      00276000
   51 CONTINUE                                                             00277000
      IF(DMAX.GT.ACC) GO TO 59                                            00278000
      WRITE(3,24) IDRAW,ICOUNT                                            00279000
      DO 60 I=1,N                                                         00280000
      P4(I)=Q1(I)                                                         00281000
      DEL(I,5)=P4(I)                                                      00282000
      WRITE(3,26) I,P4(I)                                                 00283000
   60 CONTINUE                                                            00284000
  550 CONTINUE                                                            00285000
C.....CALCULATE THE JOINT SELECTION PROBABILITIES .                       00286000
      DO 551 I=1,N                                                        00287000
      DO 552 J=1,N                                                        00288000
      IF(J.EQ.I) GO TO 552                                                00289000
      S1=0.0                                                              00290000
      S2=0.0                                                              00291000
      S3=0.0                                                              00292000
      T1=DEL(I,1)*DEL(J,2)/(1.0-DEL(I,2))                                 00293000
      IF(NS.EQ.2) GO TO 590                                               00294000
      DO 553 K=1,N                                                        00295000
      IF(K.EQ.I.OR.K.EQ.J) GO TO 553                                      00296000
      SN=DEL(I,1)*DEL(K,2)*DEL(J,3)                                       00297000
      SD=(1.0-DEL(I,2))*(1.0-DEL(I,3)-DEL(K,3))                           00298000
      T2=SN/SD                                                            00299000
      SN=DEL(K,1)*DEL(I,2)*DEL(J,3)                                       00300000
```

```
      SD=(1.0-DEL(K,2))*(1.0-DEL(K,3)-DEL(I,3))                00301000
      T3=SN/SD                                                  00302000
      S1=S1+T2+T3                                               00303000
      IF(NS.EQ.3) GO TO 553                                     00304000
      DO 554 L=1,N                                              00305000
      IF(L.EQ.I.OR.L.EQ.J.OR.L.EQ.K) GO TO 554                 00306000
      SN=DEL(I,1)*DEL(K,2)*DEL(L,3)*DEL(J,4)                   00307000
      SD=(1.0-DEL(I,2))*(1.0-DEL(I,3)-DEL(K,3))                00308000
     1  *(1.0-DEL(I,4)-DEL(K,4)-DEL(L,4))                      00309000
      T4=SN/SD                                                  00310000
      SN=DEL(K,1)*DEL(I,2)*DEL(L,3)*DEL(J,4)                   00311000
      SD=(1.0-DEL(K,2))*(1.0-DEL(K,3)-DEL(I,3))                00312000
     1  *(1.0-DEL(K,4)-DEL(I,4)-DEL(L,4))                      00313000
      T5=SN/SD                                                  00314000
      SN=DEL(K,1)*DEL(L,2)*DEL(I,3)*DEL(J,4)                   00315000
      SD=(1.0-DEL(K,2))*(1.0-DEL(K,3)-DEL(L,3))                00316000
     1  *(1.0-DEL(K,4)-DEL(L,4)-DEL(I,4))                      00317000
      T6=SN/SD                                                  00318000
      S2=S2+T4+T5+T6                                            00319000
      IF(NS.EQ.4) GO TO 554                                     00320000
      DO 555 M=1,N                                              00321000
      IF(M.EQ.I.OR.M.EQ.J.OR.M.EQ.K.OR.M.EQ.L) GO TO 555       00322000
      SN=DEL(I,1)*DEL(K,2)*DEL(L,3)*DEL(M,4)*DEL(J,5)          00323000
      SD=(1.0-DEL(I,2))*(1.0-DEL(I,3)-DEL(K,3))                00324000
     1  *(1.0-DEL(I,4)-DEL(K,4)-DEL(L,4))                      00325000
     2  *(1.0-DEL(I,5)-DEL(K,5)-DEL(L,5)-DEL(M,5))            00326000
      T7=SN/SD                                                  00327000
      SN=DEL(K,1)*DEL(I,2)*DEL(L,3)*DEL(M,4)*DEL(J,5)          00328000
      SD=(1.0-DEL(K,2))*(1.0-DEL(K,3)-DEL(I,3))                00329000
     1  *(1.0-DEL(K,4)-DEL(I,4)-DEL(L,4))                      00330000
     2  *(1.0-DEL(K,5)-DEL(I,5)-DEL(L,5)-DEL(M,5))            00331000
      T8=SN/SD                                                  00332000
      SN=DEL(K,1)*DEL(L,2)*DEL(I,3)*DEL(M,4)*DEL(J,5)          00333000
      SD=(1.0-DEL(K,2))*(1.0-DEL(K,3)-DEL(L,3))                00334000
     1  *(1.0-DEL(K,4)-DEL(L,4)-DEL(I,4))                      00335000
     2  *(1.0-DEL(K,5)-DEL(L,5)-DEL(I,5)-DEL(M,5))            00336000
      T9=SN/SD                                                  00337000
      SN=DEL(K,1)*DEL(L,2)*DEL(M,3)*DEL(I,4)*DEL(J,5)          00338000
      SD=(1.0-DEL(K,2))*(1.0-DEL(K,3)-DEL(L,3))                00339000
     1  *(1.0-DEL(K,4)-DEL(L,4)-DEL(M,4))                      00340000
     2  *(1.0-DEL(K,5)-DEL(L,5)-DEL(M,5)-DEL(I,5))            00341000
      TA=SN/SD                                                  00342000
      S3=S3+T7+T8+T9+TA                                         00343000
  555 CONTINUE                                                  00344000
  554 CONTINUE                                                  00345000
  553 CONTINUE                                                  00346000
  590 PI(I,J)=T1+S1+S2+S3                                       00347000
  552 CONTINUE                                                  00348000
  551 CONTINUE                                                  00349000
      N1=N-1                                                    00350000
```

```
      DO 556 I=1,N1                                                00351000
      J1=I+1                                                       00352000
      DO 557 J=J1,N                                                00353000
     ·PI(I,J)=PI(I,J)+PI(J,I)                                      00354000
  557 CONTINUE                                                     00355000
  556 CONTINUE                                                     00356000
      RETURN                                                       00357000
  999 IFAULT=5                                                     00358000
      WRITE(3,1000) IDRAW,MAX                                      00359000
 1000 FORMAT(1H1,////,20X,'DRAW  ',I2,'  DID NOT CONVERGE IN  ',   00360000
     1        I4,'  ITERATIONS .')                                 00361000
      RETURN ·                                                     00362000
      END                                                          00363000
      DOUBLE PRECISION FUNCTION F0(N,J,P,Q1)                       00364000
      IMPLICIT REAL*8(A-H,O-Z)                                     00365000
      DIMENSION P(N),Q1(N)                                         00366000
      F0=P(J)/(1.0-Q1(J))                                          00367000
      RETURN                                                       00368000
      END                                                          00369000
      DOUBLE PRECISION FUNCTION F1(N,J,K,P,P1,Q1)                  00370000
     ·IMPLICIT REAL*8(A-H,O-Z)                                     00371000
      DIMENSION P(N),P1(N),Q1(N)                                   00372000
      F1=P(J)*P1(K)/((1.0-P1(J))*(1.0-Q1(J)-Q1(K)))                00373000
      RETURN                                                       00374000
      END                                                          00375000
      DOUBLE PRECISION FUNCTION F2(N,J,K,L,P,P1,P2,Q1)             00376000
      IMPLICIT REAL*8(A-H,O-Z)                                     00377000
      DIMENSION P(N),P1(N),P2(N),Q1(N)                             00378000
      F2=P(J)*P1(K)*P2(L)/((1.0-P1(J))*(1.0-P2(J)-P2(K))           00379000
     1                    *(1.0-Q1(J)-Q1(K)-Q1(L)))                00380000
      RETURN                                                       00381000
      END                                                          00382000
      DOUBLE PRECISION FUNCTION F3(N,J,K,L,M,P,P1,P2,P3,Q1)        00383000
     ·IMPLICIT REAL*8(A-H,O-Z)                                     00384000
      DIMENSION P(N),P1(N),P2(N),P3(N),Q1(N)                       00385000
      F3=P(J)*P1(K)*P2(L)*P3(M)/((1.0-P1(J))*(1.0-P2(J)-P2(K))      00386000
     1                    *(1.0-P3(J)-P3(K)-P3(L))                 00387000
     2                    ·*(1.0-Q1(J)-Q1(K)-Q1(L)-Q1(M)))         00388000
      RETURN                                                       00389000
      END                                                          00390000
```

SURVEY METHOD

June/juin 1980

SURVEY RESEARCH FOR THE 1980's

## CONTENTS