

LARGE SCALE IMPUTATION OF SURVEY DATA ¹M.J. Colledge, J.H. Johnson, R. Pare, and I.G. Sande ²

Owners of small businesses complain about the quantity of forms they are required to complete and tend to blame the collectors of statistics. Administrative data are an alternative source but do not usually include all the information required by the survey takers.

The "Tax Data Imputation System" makes use of tax data collected from a large number of businesses by Revenue Canada and data obtained by sample survey for a small subset of these businesses. Survey data is imputed (estimated) for all the businesses not actually surveyed using a "hot-deck" technique, with adjustments made to ensure certain edit rules are satisfied. The results of a simulation study suggest that this procedure has reasonable statistical properties. Estimators (of means or totals) are unbiased with variances of comparable size to the corresponding ratio estimators.

1. INTRODUCTION

The demand for statistical information to aid government and management decision making has been increasing for many years. In the past, Statistics Canada was able to cope with this situation by expanding the scope and number of their surveys. Recently, such expansion has become inhibited as a result of two factors. Firstly, there is an increasing sensitivity to complaints from respondents about the burden of completing questionnaires. Secondly, current fiscal policies prevent growth in manpower. There is no indication that either of these factors is likely to be shortlived. Thus, in order to cater for an increased demand for information without raising costs or response burden, Statistics Canada is committed to making the best possible use of existing data, including data collected by other agencies for administrative purposes. One particular manifestation of this policy was the decision

¹ Adapted from a paper presented at the Annual Meeting of the American Statistical Association, August 14-17, 1978, San Diego, California, U.S.A.

² M.J. Colledge, J.H. Johnson, R. Pare, and I.G. Sande, Business Survey Methods Division, Statistics Canada.

to use financial data from Revenue Canada to supplement two annual surveys of businesses for the 1975 reference year. This paper deals with the systems which evolved as a result.

The Census of Construction (COC) is concerned with about 80,000 businesses in Canada whose primary activity is construction. The COC had been a census, but a decision was made to reduce the response burden of smaller businesses. For the 1975 reference year, only businesses with gross business income (GBI) of at least \$5,000 were considered in scope. These were divided into two groups: "small" businesses having a GBI of less than \$500,000 and "large" businesses. The latter group were the subject of a census operation; all large businesses were mailed a questionnaire asking for a comprehensive set of data. Small business information was derived from two sources: from Revenue Canada and from a mailout as follows. A sample of businesses stratified by GBI, was selected from Revenue Canada tax files, the largest business being selected with certainty. Basic financial data was transcribed for these businesses. For a subsample, secondary (more detailed) financial data was obtained from the tax return. The size of this subsample was limited by the costs of the additional transcription. A second subsample, designed to overlap the first to some extent, was selected and mailed a survey questionnaire requesting only non-financial data. The size of the second subsample was limited by the need to reduce response burden and costs. Thus in comparison with a full census, the COC response burden was reduced by sampling and reducing the number and type of questions asked.

Arrangements for the Motor Carrier Freight Survey (MCF) were along the same general lines. The significant differences were that the universe of about 25,000 was divided into "small" and "large" by a GBI threshold of \$100,000, no subsample of secondary financial data was obtained and the survey questionnaire requested a full range of information (not just non-financial).

The decision to utilize administrative tax data for the COC and MCF came quite abruptly and in advance of experience, existing software, data or feasibility study. The short time scale combined with a restricted budget dictated certain constraints on the design. Firstly, program development and testing had to be substantially achievable before any real data were available. Secondly, the programs had to be robust and easily modifiable in order to allow adjustment for unexpected characteristics of the data. Thirdly, the programs had to interface with existing systems associated with the surveys, in particular, the tabulation systems which had been developed for census operations in previous years. Thus the following design decisions were made:

- i) data from tax and survey sources would be combined at the micro level, i.e. level of individual businesses;
- ii) a complete set of data (all financial and non-financial items) would be imputed at micro level for all businesses using a "hot-deck" technique with constraints to ensure that imputation was consistent with prescribed edit rules;
- iii) the data would be inflated to universe level by replication to allow tabulation by existing systems which had not been developed to handle weights;
- iv) programs would be modular and readily adaptable to new or modified imputation and edit rules.

The following sections of this paper elaborate upon the design features and describe the systems implementation which processed 1975 data for the COC and MCF. An evaluation of the procedures is given in section 5.

2. OVERVIEW

The central feature of the system is the imputation procedure, discussed in detail in sections 3 and 4. The purpose of this section is to outline the environment within which the procedure operates by describing the complete system. The scale of processing is illustrated by reference to figures for the small business portion of the COC universe.

A system flow chart is shown in figure 1.

MERGE The first module labelled MERGE brings together data records from tax and survey sources. The input data files have been individually cleaned and edited. The output is a set of records, one per business, each of which contains a basic tax data segment and may (or may not) contain secondary tax data or survey data segments. The existing segments may have sporadic missing entries in various fields, also, some entries may be inconsistent with one another.

CHECKIN The essential purpose of the second module, CHECKIN, is to prepare data for imputation by screening out unusable or unwanted data. The module reformats the records, strips off irrelevant fields, identifies out of scope or duplicate records, checks entries against a set of prescribed edit rules, blanks out inconsistent entries and identifies all missing fields. Any record which is out of scope or a duplicate or contains insufficient useful data is flagged ("dropped"); the remainder are subject to processing by the next module, IMPUTE.

Columns 1 and 2 of figure 2 illustrate the results of processing COC data. Some 9106 of the 50,538 merged records were declared out of scope (by being in the wrong industry or too large, for example). Of the remainder, 462 were dropped leaving 40970 "good" records.

IMPUTE This is the major processing module. Its function is to impute all missing fields on every record. For the COC data, 884 records contained all segments, 3963 records required imputation of just the secondary financial segment, 2186 records required imputation of just the survey segment and 33937 records required both (see figure 2, column 3). In addition, some entries in existing segments were missing.

CHECKOUT Although, in principle, imputation is constrained by the edit rules, in practice inconsistent values may be imputed due to shortcomings in specification or programming. Furthermore, imputation may fail in the sense that no suitable value for a field can be located. Thus, the function of CHECKOUT is to check the records against the same prescribed set of rules as were applied to the data at input, and to identify and "drop" records containing inconsistent or missing entries.

From columns 2 and 3 of figure 2, it can be deduced that 194 COC records were inconsistent or incomplete and had to be dropped.

INFLATE The function of the last processing module in the system is to raise the sample of good records to the population level and thereby generate an output file which can be tabulated by the census tabulation system. Inflation is achieved by replicating each record according to its weight after "correction". All records entering the system carry a weight which is the inverse of the probability with which the record entered the basic tax sample. Three types of correction are applied prior to replication :

- i) Duplication correction. Some businesses are represented by more than one record.
- ii) Out of scope correction. There are instances where the tax data information suggests the business is in scope, whereas the survey data indicates it is not. The survey data is assumed to be more reliable. In order to allow for possible inclusion of out of scope records containing tax data only, a correction factor is applied based on data from businesses for which tax and survey information is obtained.

- iii) Dropped record correction. Records for some in scope businesses are dropped because of inadequate or inconsistent data.

Only the last type of correction is relevant in the imputation context. It implies that the imputation procedure need not be 100% successful for every record as a correction can be made.

Figure 2 indicates that after weight correction and inflation, a file of 78,563 small Construction businesses was obtained.

Imputed data is clearly identified on all files and the sponsor has access to the intermediate files to check on the reasonableness of the imputation. Some auditing and tabulation functions are also provided. The final output file has to be written in a format which can be accepted by a tabulation system which predates the imputation system and so special identifiers do not appear on this file.

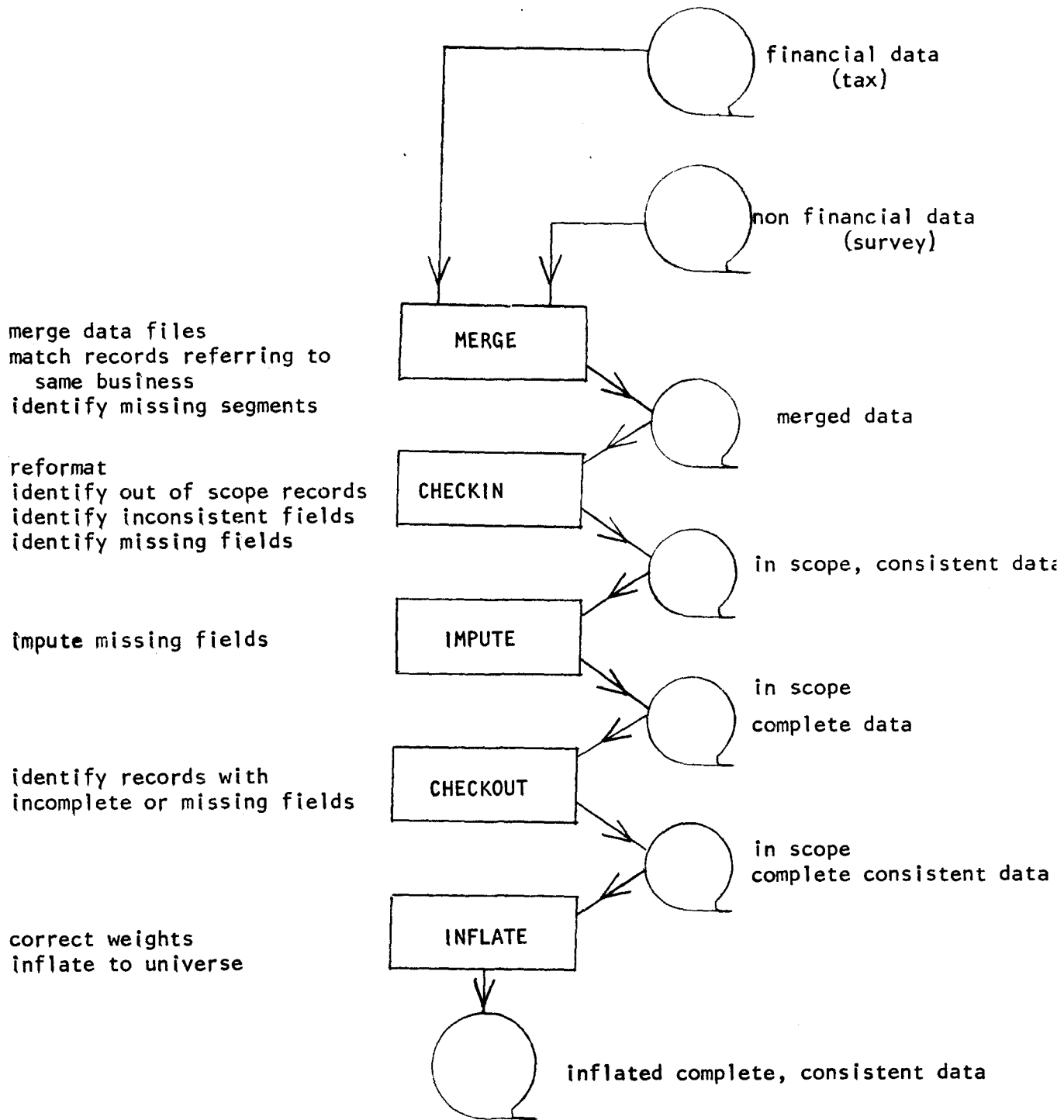


Figure 1. System Flow Chart

Figure 2. Summary of Results of Processing Census of Construction, 1975 Reference Year.

			1	2	3	4	
			At input to System	After Checkin	At output from System	Blown up to Universe	
<u>Out of Scope</u>			9106	9106	9106	-	
<u>In Scope</u>	Data not Good		0	462	656	-	
Good	Segments Present						
	Tax						
	Basic	Secondary					
	XXX		XXX	34,181	33,937	0	-
	XXX	XXX		2,316	2,186	0	-
	XXX		XXX	4,027	3,963	0	-
XXX	XXX	XXX	908	884	40,776	-	
Total Good			41,432	40,970	40,776	78,563	
Total			50,538				

3. IMPUTATION METHODOLOGY

For purposes of imputation, the record for each business can be considered as consisting of four types of segment:

- i) Key fields. These consist of fields used for classification or matching and are collected or derived from the tax return. The actual fields used were the standard industrial classification (SIC), province, salaries and wages indicator (SWI, set to 1 or 0 according as there is any indication that salaries or wages were paid or not), gross business income (GBI), net business income (NBI). If any of these fields were missing, the record was not used in the imputation.
- ii) Basic financial data collected from the tax return, e.g. depreciation, purchases, closing inventory. An attempt is made to collect this data for all businesses sampled, but the information available with the return may be insufficient or unclear. Thus some or all of these fields may be missing, i.e. the segment may be incomplete. If not, all fields are present and the segment is complete.
- iii) Secondary financial data, collected from tax returns for a subsample of records. These detailed financial data, e.g. balance sheet, detailed expense breakdowns, were collected only for the Census of Construction; but, potentially, one or more such subsamples might exist. This segment may be either complete (all fields present), incomplete (some fields present) or missing (no fields present, as in the case of records not in the subsample).
- iv) Survey data, collected for a subsample of records. This segment may be complete, incomplete or missing. In addition, there are a variety of control fields and flags.

The imputation problem is to complete the incomplete segments and to supply the missing segments.

A possible imputation procedure would be to model the missing fields in terms of those that are present. If the number of fields were very large (as it is here) and the constraints (or edit rules) on the fields were at all complex, structuring the model would be very difficult. One would have to evaluate several models to determine the best fit and this would have to be done after the data had been collected and edited. As a result, a great deal of time would be spent experimenting with the data just when one could least afford it - when the publication deadlines were approaching and a great deal of processing had yet to be done.

Thus, modelling the data did not seem a very attractive option and a type of hot-deck technique was devised. In this procedure, a record requiring imputation (candidate) is matched with a complete record (donor). The donor supplies the missing fields, possibly with some adjustment so that the edit rules are satisfied. This procedure produces realistic looking data and can be expected to preserve the underlying distributions, whereas modelling tends to produce smoothed data and distorts distributions. Another advantage is that the imputation can be set up and ready to run before the data collection is finished.

The hot-deck requires a reasonable supply of complete records, but in fact there are few records with all segments complete. If one attempted to impute for all missing fields in a single pass, the same donors would be used excessively, no use would be made of records with partial information, and the matches would be poor. In addition, a matching procedure appropriate for one segment may not be appropriate for another. Therefore, the imputation is broken up into several phases, each corresponding to a segment or sub-segment.

- Phase 1. Candidates are records with Segment A incomplete (but not missing). Donors are records with Segment A complete. At the end of Phase 1, all records have Segment A complete or missing.
- Phase 2. Candidates are records with Segment A missing. Donors are records with Segment A complete (including records which were Phase 1 candidates). At the end of Phase 2, all records have Segment A complete.
- Phase 3. Candidates are records with Segment B incomplete (but not missing). Donors are records with Segment B complete. At the end of Phase 3, all records have Segment B complete or missing. Those with Segment B complete are eligible as donors in Phase 4.

In order to match candidates with donors, the file of all records is stratified by Province (or Region), SIC and SWI. The collection of potential donors (i.e. the hot-deck) as well as the collection of candidates are identified for the particular phase. Within the stratum, the records are ordered by GBI. A sequence of records from a stratum might be represented like this:

GBI: \$25K \$26K \$27K \$28K \$29K
...CCDCCDCCD₋₅CCD₋₄CCD₋₃CCCD₋₂CCCD₋₁C₀CD₁CCD₂CCCD₃CCD₄D₅CCD

The C's are candidates and the D's are donors (other records not involved in this phase are not represented). In order to impute for C₀, only the nearest 5 potential donors on "either side" of C₀ are considered, a total of 10 possible donors which are all about the same size (in terms of GBI) as the candidate. The number 5 is quite arbitrary - it could as well be 3 or 10, or the two sides could be of different lengths, but the imputation seems relatively insensitive to this parameter. From the "nearest" 10 donors, that one is chosen which minimizes a distance function DIST (C,D). DIST can be quite a complex function, but the basic structure used was

$$\text{DIST (C,D)} = |\log \text{EXP}_C - \log \text{EXP}_D|$$

where $\text{EXP} = \text{GBI} - \text{NBI}$ = total expenses, and the subscripts C and D denote values from the candidate and donor records respectively. EXP was used because many of the fields to be imputed are detailed expense breakdowns or correlated with expenses.

Note that GBI and NBI are key fields, so that DIST is always determined. DIST may also depend on other key fields, or fields which have already been imputed in an earlier phase, or even meta-data. In particular, the distance function may be modified to spread donor usage, e.g.

$$\text{DIST (C,D)} = |\log \text{EXP}_C - \log \text{EXP}_D| (1 + p \cdot n_D)$$

where n_D = number of times the potential donor D has already been used as an actual donor in the phase,

and p = the proportional penalty for each usage (e.g. .02).

The size of p depends on the amount of imputation to be done and the degree of concern over having one donor used much more frequently than another.

After a suitable donor has been identified, the candidate's missing fields are supplied from the corresponding fields in the donor record. Some adjustment or transformation may be necessary to ensure that the constraints (edits) are satisfied. For example, three fields, X, Y and Z may have to satisfy $X + Y \leq Z$ with X, Y and Z all non-negative. The donor's values for these fields are X_D , Y_D and Z_D while the candidate has X and Y missing and the value Z_C in the Z field. If the values X_D and Y_D are simply written into the corresponding candidate fields, we may find that $X_D + Y_D > Z_C$, which violates the edit. Therefore, it is better to prorate X_D and Y_D to ensure that the edit holds:

$$X_C = (X_D/Z_D) Z_C; \quad Y_C = (Y_D/Z_D) Z_C \quad .$$

In other words, the proportions X_D/Z_D and Y_D/Z_D are transferred to the candidate. A common example is

$$FUEL_C = (FUEL_D/EXP_D) EXP_C$$

where FUEL is the amount spent on fuel and EXP is the total expenses. This imputation estimates that the candidate spent the same proportion of his total expenses on fuel as did the donor.

The transformation needed to impute a field may be more complex if the field is involved in several edits. For example, the four fields W, X, Y, Z, may have to satisfy $X + Y \leq Z$ and $X \leq W$, where all fields are non-negative. The donor's values for these fields are W_D, X_D, Y_D, Z_D . The candidate has W_C, X and Y missing, and Z_C . An appropriate imputation (but not necessarily the only one) is

$$X_C = \text{Min } W_C, (X_D/Z_D) Z_C$$

$$Y_C = (Y_D/Z_D) Z_C.$$

When the edit rules are even more complex a decision table may be required, where the form of imputation depends on which set of conditions holds. In desperate situations, a table of default values may be used.

If a field is not involved in any edits, it may be prorated using a correlated variable in the case of a numeric field. Categorical data may simply be copied from donor to candidate.

The imputation specifications are written separately for each field - no generalized transformation is used. They are written in such a way as to produce consistent data and this involves not only accommodating constraints, but also ensuring that constraints are not violated due to roundoff error.

4. IMPLEMENTATION

The systems design was based on the following premises:

- a) The breakdown into phases each of which is functionally the same, except in detail, suggested a general system which would be tailored separately for each phase.
- b) To simplify data-set control, the output produced from a phase would have the same record description as the input and all records would be carried forward. Each phase would identify its donors and candidates, perform imputation, and copy all other data as is.
- c) Instrumentation of the system would mostly be done offline by analysis of a log file describing imputation "events", and by investigation of the output of each phase.
- d) Fields would either have a value or be missing. If missing, any value which it might have had would be ignored for imputation purposes.
- e) Fields would be identified as missing only at beginning of processing. Once imputed to a value, the field stays imputed. Thus, inconsistencies must be removed at the beginning and never introduced by imputation.
- f) The control language should be quite flexible to allow unusual imputation rules, but should still be quite readable since it would be the final specification of side effects in unusual situations.
- g) One donor only would be used in each phase.

The effect of these considerations on the design was to simplify the systems development and operation of the system while retaining flexibility in the details of imputation. This would facilitate final turning without holding up production more than necessary.

Consideration a) resulted in the general phase structure shown in Fig.3. Basically four modules are involved along with three utility sorts:

- i) CNVT is responsible for identifying that subset of the file that is to be involved in imputation. For each donor or candidate it writes out an "Imputation Control Segment" (ICS) which contains match fields for donor assignment as well as space for indicating the donor actually assigned.

- ii) NEBR performs the assignment of donor to candidate on the basis of match fields. The ICS file has been stratified by sorting on a KEY. A local search is performed in a large circular buffer (about 2000 segments) and the best match according to some measure is selected.
- iii) MERG combines a copy of the appropriate donor record to each ICS record.
- iv) IMPT then performs consistent imputation using the donors assigned.

Consistent imputation (for linear edits) was aided by a routine that kept track of the current upper and lower bounds for each field, determined by the edits and the fields already assigned. For each field to be imputed, assignment would be done if the value were in range, and the ranges of the remaining unassigned fields would then be adjusted appropriately. The routine caused the actual assignment to be made and a log entry to be written.

Where it could be applied, this approach simplified the work enormously. Unfortunately, it could not be made universally applicable without in effect solving an integer programme at each field assignment. Nonetheless, the edit rules which occurred were predominantly positivity restrictions and simple sums. Some conditional edits could be handled by selectively activating edits. Others were handled by taking great care with the imputation rules. However, the potential for an inconsistent imputation still remained.

Flexibility (consideration (f)) was ensured by allowing the control language to be a number of inclusions into the general programmes which could then be compiled to produce executable modules. The environment of each inclusion is carefully documented and service routines are provided for certain common functions.

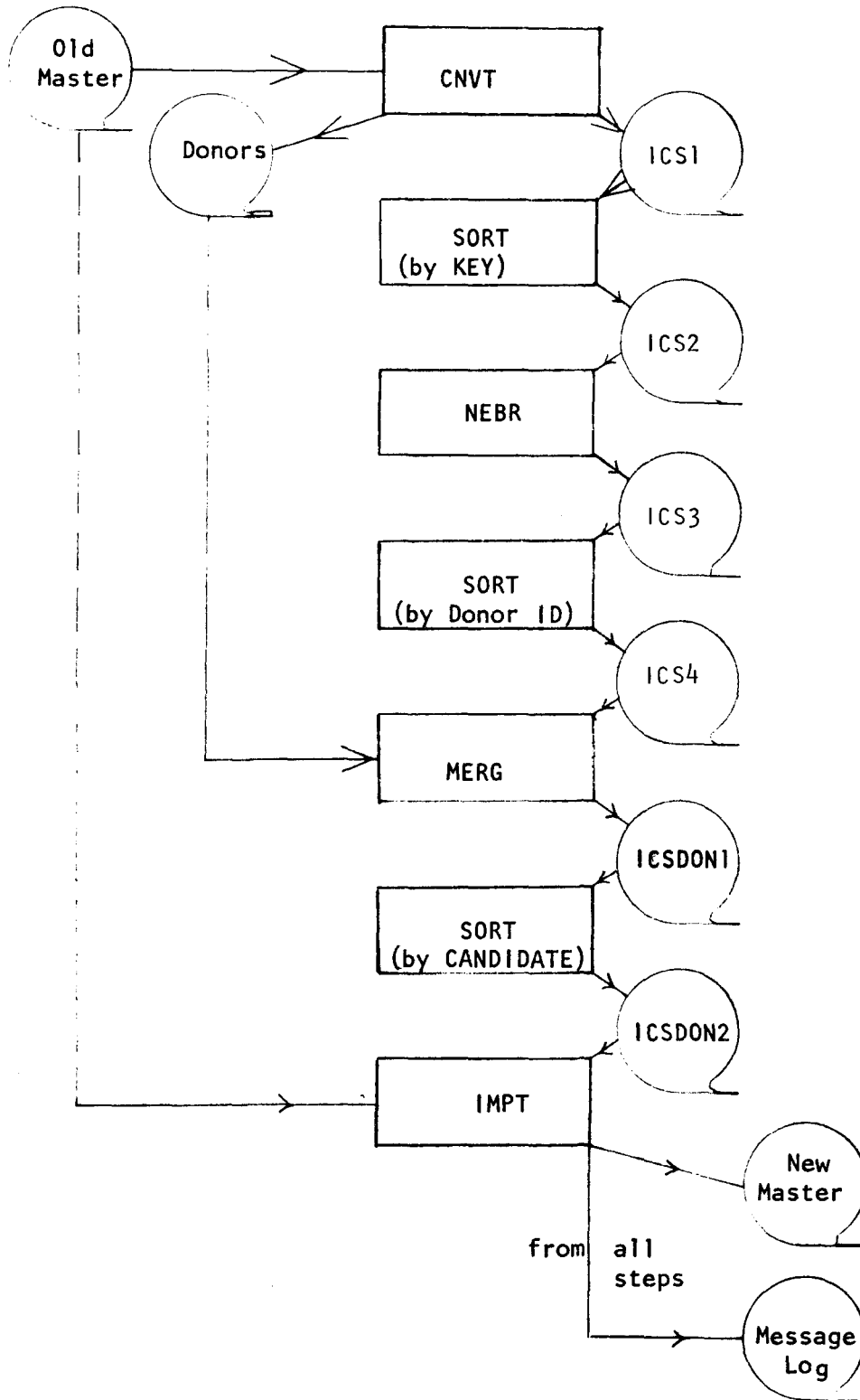


Figure 3. General Phase Structure

5. EVALUATION

The imputation procedure described in section 3 will produce estimates of the population totals (or means), but some assessment of the quality of these estimates, in terms of bias and variation, is required. One would like to know how the quality of the estimate varies with (a) the sampling bias, (b) the population size, (c) the sampling rate, (d) the correlation or relationship between the imputed variable and the auxiliary variable used for prorating, (e) the size of the window used to determine the number of eligible donors, (f) the complexity of the edits, (g) the distance function, and (h) the control of donor usage. One would also like to compare the "imputation" estimate with some natural competitors, such as the usual sampling (expansion) estimate and the ratio estimate.

A small simulation study has been done to examine the effects of sampling bias (in a nominally simple random sample) and sampling rate for a population of fixed size.

A population of 1000 units was created, each consisting of five variables corresponding to GBI, NBI and the "expense items": "salaries", "depreciation" and "purchases". GBI and NBI were the auxiliary variables. All quantities except NBI are non-negative and, in addition, we have the edit rule.

$$\text{Salaries} + \text{depreciation} + \text{purchases} \leq \text{GBI} - \text{NBI} = \text{EXP.}$$

We omit the gory details, but the distribution of the non-negative variables is skewed towards zero.

Sampling was either unbiased or biased. Biased samples were created by ordering the population on GBI and (a) selecting 25% of the sample from below the median GBI and 75% of the sample from above the median GBI (bias up), or (b) reversing the percentages in (a) (bias down).

The sampling fractions were 10%, 20% and 50%.

For each sampling bias and sampling rate, twenty-five independent samples were selected from the same population. For each sample, a new file was created for the population in which GBI and NBI were retained for all records and salaries, depreciation and purchases were included for the sampled records only. Salaries, depreciation and purchases were then imputed for the non-sampled records, using the sampled records as the hot-deck and prorating on EXP. For each replicate, the imputation, sampling and ratio estimates of the population means were calculated. These could then be compared with the known population values.

Table I gives the mean over 25 replicates divided by the population mean for each type of estimate, bias condition, sampling rate and variable. The t statistic, evaluating the "significance" of the difference between the population mean and the average value of the 25 estimates is given in parenthesis. The population correlation between the imputed variable and the prorating variable is given in parenthesis in the first column. For the unbiased case, all types of estimates do quite well, except that the ratio estimate begins to show bias at a 50% sampling rate. For the biased cases, the imputation estimate clearly does better than the ratio estimate. The sampling estimate does very badly as one would expect.

Table II gives the coefficient of variation of the estimates in the form of the standard deviation calculated for the 25 replicates divided by the population mean. For the unbiased case, the coefficients of variation are about the same for the imputation and ratio estimates, while that of the sampling estimate is much larger. This is also true for the upward biased case. In the downward biased case, the position is less clear and the estimates appear to be roughly equivalent; but if one considers the root mean square error divided by the population mean, the bias dominates and the imputation estimate is clearly superior.

The implication of Table II is that in order to estimate the variance of an imputation estimate (in a "real" situation where replicates are not available) one may formally use the estimate of the variance of the corresponding ratio estimate as a reasonable approximation.

It will be noticed in Table I that the correlations between the imputed and prorating variables are quite high, higher than one might expect in "real" data. We would expect the difference between the imputation and the ratio estimate to become less pronounced as the correlation decreased; but no systematic work has been done to investigate this.

When the correlations are high, the size of the window appears to have no effect on the quality of the imputation estimate.

We have some evidence to suggest that when the correlations are low and the sampling rates are very low, all estimates are bad.

Table 1 : Relative Bias of Estimates

Mean over 25 replicates / Population mean (t 24)

Variable ρ	Sampling Fraction	UNBIASED			BIASED UP			BIASED DOWN		
		Imput- ation	Sampling Ratio	Ratio	Imput- ation	Sampling Ratio	Ratio	Imput- ation	Sampling Ratio	Ratio
Salaries (.95)	10%	.996	1.004	.998	.995	1.329	1.026	1.005	.690	.961
		(-.6)	(.2)	(-.2)	(-1.1)	(16.6)	(4.2)	(.6)	(-23.6)	(-4.8)
	20%	.998	1.000	.999	.997	1.330	1.029	.999	.677	.947
		(-.4)	(0.0)	(-.2)	(-.7)	(32.0)	(8.9)	(-.1)	(-34.7)	(-6.9)
	50%	1.000	1.000	.996	.996	1.338	1.030	.996	.670	.944
		(-.1)	(-.02)	(-3.4)	(-2.1)	(55.7)	(21.8)	(-1.2)	(-115.4)	(23.9)
Depreciation (.89)	10%	1.004	1.003	1.000	1.004	1.254	.993	.993	.746	1.041
		(.5)	(.2)	(.1)	(.7)	(30.2)	(-9.0)	(-.8)	(-25.5)	(4.9)
	20%	1.001	1.000	1.001	1.004	1.251	.968	1.002	.754	1.056
		(.2)	(0.0)	(.2)	(.8)	(43.8)	(-9.0)	(.3)	(-47.4)	(7.0)
	50%	1.000	.993	1.003	1.005	1.258	.969	1.004	.751	1.059
		(-.2)	(-1.4)	(2.0)	(2.0)	(88.1)	(-17.2)	(1.1)	(-80.7)	(23.8)
Purchases (.82)	10%	.993	1.008	1.002	1.000	1.342	1.036	1.004	.681	.946
		(-.6)	(.3)	(.2)	(0.0)	(12.9)	(2.8)	(.2)	(-17.3)	(-3.2)
	20%	.999	1.000	.999	.993	1.341	1.038	.979	.659	.921
		(-.1)	(0.0)	(-.1)	(-.9)	(23.4)	(5.1)	(-1.4)	(-25.9)	(-6.2)
	50%	.999	.996	.992	.995	1.343	1.034	.994	.658	.927
		(-.3)	(-.6)	(-2.5)	(1.8)	(45.1)	(12.7)	(-.8)	(-72.4)	(-12.7)

Table 11 : Coefficients of Variation of Estimates

Standard Deviation of 25 Replicates/Population mean.

Variable	Sampling Fraction	Imputation	UNBIASED Sampling Ratio	Imputation	BIASED UP Sampling Ratio	Imputation	BIASED DOWN Sampling Ratio			
Salaries	10%	.039	.106	.039	.024	.099	.031	.043	.066	.041
	20%	.021	.081	.022	.023	.052	.017	.040	.047	.038
	50%	.010	.028	.006	.011	.030	.007	.018	.014	.012
Depreciation	10%	.039	.078	.038	.024	.042	.031	.043	.050	.041
	20%	.021	.059	.021	.024	.029	.018	.040	.026	.040
	50%	.010	.025	.007	.012	.015	.009	.019	.015	.012
Purchases	10%	.066	.127	.065	.055	.132	.065	.101	.092	.084
	20%	.039	.091	.042	.039	.073	.037	.073	.066	.065
	50%	.021	.036	.016	.014	.038	.013	.041	.024	.029

6. CONCLUSION

Planning for the 1975 imputation system started in April 1976 and the final output data were delivered in August 1977. Most of the delays were due to problems with data collection and survey processing. Publications based partly on the imputed data have been released.

For 1976 data, the imputation system and methodology were refined and at least one survey, the Census of Construction, should run on virtually the same system with 1977 data.

Large-scale imputation appears to be a useful new weapon in the arsenal; but more evaluation should precede more widespread use. At the moment, assessment of its feasibility in any situation is based more on hunches than facts. Unfortunately, thorough and systematic evaluation promises to be a lengthy process and the best we can hope for are piecemeal results.

RESUME

Les petits entrepreneurs se plaignent de la quantité de formules qu'il leur faut remplir et ont tendance à accuser les responsables de la collecte des statistiques. Les dossiers administratifs constituent une autre source possible, mais il y manque souvent des renseignements essentiels aux enquêteurs.

Le système d'imputation à l'aide des données fiscales a recours aux données fiscales recueillies par Revenu Canada auprès d'un grand nombre d'entreprises et aux données obtenues par sondage auprès d'un petit sous-ensemble de ces entreprises. Les données sur les entreprises qui ne font pas partie de l'échantillon du sondage sont imputées (estimées) par la méthode du hot-deck, certaines corrections étant apportées pour assurer le respect de diverses règles de validation. Les résultats d'une simulation semblent indiquer que cette méthode possède des propriétés statistiques raisonnables. Les estimateurs (des moyennes ou des totaux) sont sans biais, et leurs variances présentent des grandeurs comparables à celles des variances des estimateurs obtenus par la méthode du quotient.