# CUT-OFF SAMPLING AND ESTIMATION

Hanna Elisson[1] and Eva Elvers[2]

## ABSTRACT

In cut-off sampling, part of the target population is deliberately excluded from selection; in business statistics the frame and the sample are typically restricted to enterprises of at least a given size, e.g. a certain number of employees. The response burden is eliminated for the small enterprises, but assumptions must be used for the non-sampled part of the population. Some empirical Swedish results are presented based on one survey and administrative data. Different error sources and their effects on the overall accuracy are discussed. Cut-off sampling has merits but requires care when measuring size and methodological work with models.

KEY WORDS:        Cut-off sampling; Administrative data; Model-based estimation

## 1. INTRODUCTION

In business statistics it is not unusual to cut off (very) small enterprises from the sampling frame. The contribution from this part of the population is, if not negligible, at least small in comparison with the remaining population. It may be tempting not to use resources on enterprises that contribute little to the overall results of the survey. Moreover, this reduces the response burden for these small enterprises. On the other hand, cut-off sampling may be considered a dirty method, simply because the sampling probability is set equal to zero. Should we consider cut-off sampling as a permitted method or not? If we do, how should the cut-off threshold be chosen?

When trying to answer these questions about cut-off sampling and appropriate thresholds, we have not found much experience and practice to build on. There is a description of cut-off sampling in the book by Särndal, Swensson, and Wretman (1992) with some warnings. Haan, Opperdoes, and Schut (1999) have positive experience with cut-off sampling in a Consumer Price Index context. The Canadian Monthly Survey of Manufacturing (MSM) has made a choice: "The sampling frame for the MSM is determined from the target population after subtracting establishments that represent the bottom 2% of the total manufacturing shipments estimate for each province. These establishments were excluded from the frame so that the sample size could be reduced without significantly affecting quality." The citation is taken from the documentation series Data Quality Statements, Statistics Canada (2001).

At Statistics Sweden several surveys use cut-off sampling. The technique has a tradition, but we have not found any methodological documentation behind it. Lately the threshold has been increased in a few surveys as a way to manage reduced budgets. This was done at a short notice, too short in our view. Some more surveys are now under pressure. The topic of cut-off sampling deserves methodological studies, which consider different advantages and disadvantages with regard to quality and costs. There should be clear principles when using cut-off sampling and for determining the threshold.

Here we present the first findings in our work with cut-off sampling. We have chosen one Swedish survey for numerical illustrations. As is often the case in a methodological study, a study of one problem leads to additional problems and findings. We start by describing some Swedish data in Section 2 in order to give

[1]        Hanna Elisson, Statistics Sweden, SCB ES/SES, Box 24 300, SE-104 51 Stockholm, Sweden
[2]        Eva Elvers, Statistics Sweden, SCB ES/SES, Box 24 300, SE-104 51 Stockholm, Sweden

essential background information to our studies. Next, in Section 3, we discuss accuracy from the point-of-view of cut-off sampling. Section 4 deals with fairly small enterprises and Section 5 with those, which have no or few employees. In Section 6 Neyman allocation is used to test the importance of different size-classes and to compare variances. Section 7 concludes.

## 2. BASIC SWEDISH INFORMATION

### 2.1 Administrative data, the Business Register, and the sampling system

The Swedish Business Register (BR) is a register of several statistical units, notably enterprise, local unit, kind-of-activity unit (KAU), and local kind-of-activity unit. An enterprise normally consists of one legal unit, but in certain cases an enterprise has more than one legal unit. There are now around 50 composite enterprises with altogether about 600 legal units in the BR. The present Swedish practice is a bit different from the recommendation of the European Union. Even if the composite enterprises are small in number, they are mostly important in their industry. In this study of cut-off sampling the interest focuses on small and fairly small enterprises, so we will concentrate on enterprises with just one legal unit and one KAU.

There is information about turnover and employment for legal units, through VAT (value-added tax) and PAYE (administrative information from the collection of taxes on earnings, which includes employment; Tax Payroll). A legal unit has an identification number used by the fiscal authorities and in the BR. The BR obtains information about births and deaths of legal units from the National Tax Board every second week. The number of employees is updated through several sources. The two main ones are PAYE-information, which is used to compute the number of employees for single-location legal units, and a BR questionnaire to multiple-location legal units. Each of these is essentially performed once a year. Statistics Sweden gets VAT information monthly. The legal unit is basic and indeed a valuable source of information. However, different legal units can report VAT and PAYE for one and the same activity. This is the case for legal units within a composite enterprise, but also for other groups of legal units.

Statistics Sweden has a sampling system for business statistics that most regular surveys use to draw samples. The system uses so-called permanent random numbers (PRN's). This is a convenient and flexible method of coordination. There is both positive and negative co-ordination between surveys. Each survey gets a high overlap over time between its successive samples, which is favourable when estimating changes. There is a rotation system, though, in order to reduce the response burden.

Many samples to be used for short-term statistics in year t are drawn in November year (t–1). The information of the frame is fairly recent for some variables: the end of September for active enterprises and local units. Other information is older: the number of employees refers to the spring in year (t–1) for multiple-location enterprises and to the end of year (t–2) for single-location enterprises (BR questionnaires and PAYE information, respectively). From 2001 and onwards, frames are created also in March. The main advantage is the inclusion of a considerable number of reorganisations by January 1[st]. Single-location enterprises normally have 0 employees in the BR in the year of birth until May next year. This delay contributes to the heterogeneity in the size-class with zero employees. Surveys that require a minimum of for example 10 employees normally do not cover births in the previous year.

The methodologists highly recommend more frequent sample updates than once a year. There is an increasing understanding and tendency to follow the advice. Some short-term surveys draw samples twice a year. Since some important variables in the BR are updated only once or a few times a year, there is not much to gain from very frequent renewals of the sample.

## 2.2 The survey used for numerical illustrations

The survey that we use as an example in our computations has been through a number of changes during the last five years. There used to be three separate surveys: new orders and deliveries (monthly), inventories (quarterly), and capacity utilisation (quarterly). These three surveys were integrated into one survey in 1998. There is a monthly questionnaire and an additional questionnaire every third month. One of the aims of the integration was to enhance production (as opposed to deliveries). The monthly questionnaire includes the number of production days and the production value, also the values of deliveries, new orders, and the stock of orders (split into domestic market and exports). The "extra" quarterly questionnaire includes values of inventories of different types at the end of the quarter and changes during the period. The questionnaire layout underlines the relationships between the variables – to assist and to make measurement errors small.

The statistics cover mining and quarrying and manufacturing. The industrial classification used by Sweden is based on the European NACE Rev. 1 (Statistical Classification of Economic Activities in the European Community). The NACE code has two letters and four digits, and the Swedish version has an additional fifth digit in some cases. The observation unit is roughly equal to the kind-of-activity unit (the difference has historical reasons and is disappearing). The variables, such as value of deliveries, refer to own manufacturing, not trade. The level of detail of the statistics is determined by users' needs, especially those expressed by the National Accounts and Eurostat. The survey design depends on these requests and data collection possibilities. Stratified simple random sampling is used. 49 industries are crossed with 6 size-classes, based on the number of employees. Table 1 shows all size-classes, including those cut-off.

Table 1. Size classes and sampling

| Size-class number | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| Number of employees | 0 | 1–4 | 5–9 | 10–19 | 20–49 | 50–99 | 100–199 | 200–499 | 500– |
| Sample | Cut-off | | | Take-some | | | | Take-all | |

The non-response rate increased considerably in the late 1990's. Several surveys had difficulties at that time to maintain the response rate. For this particular survey, the change into new questionnaires was an additional cause. Especially the production variable has considerable item non-response. For many respondents the information required is not part of the ordinary accounting system. Some efforts have been made to decrease the response burden, mainly by reducing the number of variables for the smallest enterprises, which get a simplified questionnaire. There is a general pressure to reduce the response burden, especially for small enterprises. There is also pressure both to decrease the working load on the survey staff and to improve quality. As a way to save resources the survey staff has suggested a higher cut-off threshold. There is both high unit non-response and high item non-response in the size-class 10-19 employees.

We work with survey data and administrative data. There are no survey data for the smallest size-classes and quite few observations just above the present cut-off threshold. A comparison between the turnover from VAT and the survey variable "deliveries from own manufacturing" shows a number of facts. The figures are very close in most cases. The turnover is a bit higher, especially for large and composite enterprises, as expected from the definitions. For a small proportion of enterprises the difference between the values is substantial. This is likely to be due to different uses of reporting units.

## 3. ACCURACY AND SOURCES OF INACCURACY FOR CUT-OFF SAMPLING

## 3.1 Introduction

When considering accuracy, we distinguish between systematic and random errors. We strive for procedures without bias and with a small variance. The mean squared error is a way to summarise the two types of inaccuracy. When analysing and describing accuracy it is convenient to work with the sources of

inaccuracy. There are six such sources shortly named as follows: sampling, frame coverage, measurement, non-response, data processing, and model assumptions. The contributions to inaccuracy from these sources depend on many factors; one such factor is the sizes of the enterprises in the sample. The largest enterprises can be disregarded in our discussion related to cut-off sampling, since they are always included.

We compare two main possibilities for a survey. The first possibility is a sample survey for the whole population, using stratified random sampling, with the stratification based on industry and size-class. The second possibility is a similar sample survey but restricted to the population above a cut-off threshold and complemented with model-based estimation for the cut-off part. The survey has a budget, and the idea is to achieve as high accuracy as possible given that budget. We are interested in effects of introducing cut-off sampling and in further effects when increasing the cut-off threshold.

## 3.2 Inaccuracy sources and size-dependence

The frame population gets smaller if cut-off sampling is introduced, and it decreases with an increasing cut-off threshold. If the sample size is the same, the inaccuracy due to *sampling* also decreases. The sample size is not necessarily exactly the same, but with a given budget the variation in sample size is likely to be small.

Deficiencies in *frame coverage* – differences between the frame population and the target population – are highly due to delays in receiving information. There are also pure, remaining errors. The population of enterprises changes quickly: births, deaths, reorganisations, and changes in activities and size. Over- and under-coverage depend on size, and the proportions may be fairly high for small enterprises. When cut-off sampling is introduced – for example a minimum of z employees – the coverage of the surveyed population is to be measured with this restriction. Enterprises in the survey that have less than z employees are over-coverage, and enterprises outside the survey with at least z employees are under-coverage. In a previous study, Elvers (1993), this error source was found to be substantial for a Swedish investment survey with the cut-off threshold at 20 employees. Observed stratum mean values were too low in the smallest size-classes. These size-classes accounted for a moderate share of the investments, but the effect of coverage deficiencies was serious. The time schedule for taking frames has been changed.

When cut-off sampling is used, a simple and crude way to proceed is to neglect the part of the population that was cut off. Even though not necessarily stated explicitly, this is a type of *model assumptions*. When estimating a total, the implied assumption is "negligible contribution". When estimating an index of change, the implied assumption is: "the growth rate is the same below and above the threshold". In both cases an alternative way to proceed is to build an explicit model where survey and register information are utilised. It should be possible to improve the accuracy in this way. The accuracy can be estimated through the model itself and/or by using external information. There may be improved estimates and evaluations later on, when more information is available.

The size-dependence for the inaccuracy components *non-response*, *measurement*, and *data processing* can be seen in a few different ways. The final response-rate often increases with enterprise size, but this is partly due to the strategy of reminders, so the cost is also size-dependent. There are difficulties to get responses and adequate responses, for example in cases of bankruptcy, when there is some but not full activity. Again, there is a dependence on size for both quality achieved and costs for reminders, follow-ups etc. We have few hard facts about these quality components, and even less for the size-dependence.

The above reasoning shows that both the quality achieved and the cost depend on the sizes of the sampled enterprises, although our knowledge is limited. In our first step we concentrate on the sampling error and on the model needed when using cut-off sampling. Also as a first approximation, we will consider the overall sample size to be determined by the budget. Hence, we take the sample size to be the same both without and with cut-off sampling. As already stated, there will be differences at a closer look depending on need for re-contacts, reminders, checks about deaths and other causes of over-coverage etc. It seems reasonable that there are differences between size-classes, some considerable enough to take into account.

# 4. TURNOVER PER EMPLOYEE FOR FAIRLY SMALL ENTERPRISES

As mentioned above, the survey staff considers increasing the cut-off threshold from 10 to 20 employees. One question is if it is possible to make a model-based estimate of the total turnover for enterprises with 10-19 employees using information about enterprises with 20-49 employees. A possible and convenient assumption is that turnover per employee is the same within industry. There is only a small amount of data collected in these size-bands; therefore we have mainly used administrative data (turnover from VAT) in our calculations. We used two different models to estimate the total turnover in size-class 3.

Model 1.
The total turnover in size-class 3 is estimated by multiplying the number of employees in this size-class by the estimated turnover per employee in size-class 4. Simple regression is used with annual turnover as the dependent variable y and the number of employees as the independent variable x. The model is without intercept, and the residual variance is proportional to x. The regression parameter, $b_{jg}$, is estimated per group, $U_{jg}$, where j denotes size-class and g denotes industry-group. An industry-group consists of at least one stratum. The model-based estimator of the total Y for size-class 3 and group g is given by

$$\hat{Y}_{3g}^{mod\,1} = \hat{b}_{4g}^{t} \sum_{k \in U_{3g}} x_k \qquad \text{where} \qquad \hat{b}_{4g}^{t} = \sum_{h \in U_{4g}} \sum_{k \in s_h} \frac{N_h}{n_h} y_k \Bigg/ \sum_{h \in U_{4g}} \sum_{k \in s_h} \frac{N_h}{n_h} x_k$$

and $s_h$ denotes the sample selected from stratum h, $N_h$ is the population-size and $n_h$ is the sample size. Time is not shown explicitly, except for the year t, needed in model 2.

Model 2.
A closer analysis of data for three successive years showed that enterprises are "individuals". There are differences not only between industries but also within industries. Some enterprises have fairly high values of "turnover per employee", while others have fairly low values. Enterprises with high (low) values one year often have high (low) values the following years. Hence the ratio $Q=b_{3g}/b_{4g}$ between the parameters in size-classes 3 and 4 is fairly stable over time for the same set of enterprises. We can estimate $Q$ by using turnover from VAT for a previous period. We can apply this $Q$ in the current period where we have collected data for size-class 4 but not 3. We do so for enterprises that are in the sample on both occasions; there are fairly many such enterprises due to sampling with positive coordination. In comparison with the first model an "extra" factor $q_k$ is introduced, which deviates from 1 and equals Q for these enterprises. The second model-based estimator of the total Y for size-class 3 and group g is given by

$$\hat{Y}_{3g}^{mod\,2} = \hat{b}_{4g}^{t} \sum_{k \in U_{3g}} x_k\, q_k \qquad \text{where} \qquad q_k = \hat{b}_{3g}^{t-1} \Big/ \hat{b}_{4g}^{t-1} \quad \text{if} \quad k \in U^{t-1} \qquad \text{and } q_k = 1 \text{ otherwise,}$$

and $U^{t-1}$ is the population on occasion *(t-1)*.

These two models have been used to estimate the total turnover in size-class 3. Some results are presented in Table 2. The two model-based estimates are compared with the known total turnover. The comparison is performed on two levels: size-class 3 and size-classes 3-8, i.e. enterprises with 10-19 employees and at least 10 employees, respectively. The five industries have been selected to show a variation of results. The proportion of employees in size-class 3 is more than 17 percent of the total number of employees for enterprises in the three industries NACE 18, 19, and 28. It is only one percent for NACE 21 and 34 – these industries are not very sensitive to the choice of model. Industries NACE 18 and 19 have very low totals and are from that point of view less important than NACE 28.

Model 2 is mostly better than model 1. The deviations are very small for NACE 34, not only for the "whole" industry, but also for size-class 3. Model 2 works well also for NACE 28, which is in the group with a high share of small enterprises. It should be observed that the present size of inaccuracy due to sampling is at least ten times the difference between the known total and the model-based estimate.

Table 2. Numeric results using model-based estimates.

| NACE | | | Enterprises with 10-19 employees | | | Enterprises with at least 10 employees | | |
|---|---|---|---|---|---|---|---|---|
| | | | Total turnover | Estimated total turnover | Absolute difference in percent | Total turnover | Estimated total turnover | Absolute difference in percent |
| 18 | Manufacture of textiles and textiles products. | Model 1 | 401 | 390 | 0.03 | 1781 | 1770 | 0.01 |
| | | Model 2 | | 327 | 0.18 | | 1707 | 0.04 |
| 19 | Tanning and dressing of leather. | Model 1 | 233 | 129 | 0.45 | 1167 | 1063 | 0.09 |
| | | Model 2 | | 136 | 0.42 | | 1070 | 0.08 |
| 21 | Manufacture of pulp, paper and paper products. | Model 1 | 1126 | 1201 | 0.06 | 102453 | 102528 | 0.0007 |
| | | Model 2 | | 1208 | 0.07 | | 102536 | 0.0008 |
| 28 | Manufacture of fabricated metal products. | Model 1 | 10059 | 10890 | 0.08 | 71057 | 71888 | 0.01 |
| | | Model 2 | | 10255 | 0.02 | | 71253 | 0.003 |
| 34 | Manufacture of motor vehicles, trailers and semi-trailers. | Model 1 | 1023 | 940 | 0.08 | 193455 | 193372 | 0.0004 |
| | | Model 2 | | 1024 | 0.001 | | 193456 | 0.00001 |

## 5. DIFFICULTIES ESPECIALLY WITH SMALL ENTERPRISES

When building an estimation model we aim for characteristics, which can be used across size and which are based on information available in registers. When forming the ratio between turnover and number of employees, we use data that are easily available, but we are aware of the fact that the denominator is not appropriate. We would like to have a measure of the work performed, and the number of employees is a substitute with several deficiencies. There may be further working persons who are not employees – for example sole-proprietors – and the amount of work per person varies, since a person can work full-time, part-time, temporarily etc. There are variables in the BR, e.g. legal form, which provide some information. We need to modify the number of employees – especially when the number is zero – to get a number more similar to the number of working persons. We now simply add a number to the number of employees: a number, which equals one for enterprises without employees and which decreases with the number of employees down to zero for ten employees.

The VAT and PAYE data sets refer to different time periods. If the legal unit has undergone changes (merges, splits etc.), the ratio between the turnover and the (modified) number of employees is not meaningful. The match of identification numbers is only formal, since the values in the numerator and the denominator mismatch. Some such formal matches are temporary, while others last for many years. Groups of legal units choose to report turnover to fiscal authorities where suitable for tax reasons. The enterprise is better for matching than the legal unit, but there are still many mismatches of activities. We were aware of the problems with mismatches before we started our work, but both the number of units and the values involved are higher than we had expected. One of our first findings when running the allocation program described in Section 6 was the effect of extreme values. There are turnover values corresponding to 500 and 50 employees among enterprises without employees. Although they are not very many, this is a problem that influences our cut-off study. It has other implications as well, for several surveys.

An analysis is needed of the extreme values and the legal units involved. Some of the high turnover values (but not the extreme ones) have been in the size-class without employees for some time. The most extreme values are explained by recent reorganisations. We believe that most surveys have detected most of these changes, but there may be delays. One of the aims when introducing a new BR in 2000 was harmonisation. The information flow between the BR and the surveys has certainly been improved, but more can be done. We also draw the conclusion that methodologists should examine the survey design. VAT-turnover should perhaps be used. A simple possibility is to add enterprises with high turnover and few employees as a special group. If so, care is needed to avoid double counting. Alternatively, turnover could be the only, or dominating, measure of size. A cautious study is needed before a change; considering how and when these variables are updated in the BR and how they relate to survey variables.

In Section 5, an estimation model for enterprises with 10-19 employees was derived. Before working on a model for smaller enterprises, the problem with extreme and high values has to be handled. The small size-classes are heterogeneous, especially the class without employees. There is a small number with high turnover values. Moreover, there is a high number with low values. The simple ratio between the turnover and the modified number of employees taken per size-class increases for the size-classes, within industry.

## 6. A TEST WITH NEYMAN ALLOCATION

We have created the "full" frame for our illustrative example. The frame is from May 2001, and it has values of turnover from VAT for the year 2000. We have deleted enterprises without both VAT and employees from our computations. Otherwise we have replaced missing values for VAT with nil turnover value. There are then about 31 000 enterprises in mining and quarrying and manufacturing. We have 49 industries crossed with the 9 standard size-classes shown in Section 2.2. Enterprises with 0-9 employees are below the present cut-off threshold. They account for about 5 % of the turnover and 75 % of the enterprises.

As a simple test of the importance of different size-classes, we have used Neyman allocation for a given sample size. We use VAT-turnover for this allocation, and we "remove" the problem of extreme values described in Section 5 by creating an extra size-class, where we sample all units. This is a simple way to keep the turnover values by industry without worrying about their exact and adequate location (statistical unit and size-class). As stated in Section 5, this group of enterprises needs an investigation. The turnover variable is close to one of the survey variables. According to the Neyman allocation rule, each stratum gets a sample size, which is proportional to the product of the stratum population size and the stratum population standard deviation of the allocation variable (see e.g. Särndal, Swensson, and Wretman, p. 106).

We aim for a total sample size that is similar to the present one. We do not yet work with costs. We have so far used precision requests for industries on the two-digit level of the NACE code. We have formulated a rule, and we have run the allocation program without cut-off and with a few different cut-off thresholds according to the standard size-classes. The total sample size is the same each time, 2 200 enterprises. This figure does not include the extreme values, but it is lower than the present sample size to have a safety margin. We make the comparisons step-wise.

The first step is to cut off the smallest size-class, the class without employees. There is then a clear decrease in sampling variance for most industries on the two-digit level. We can compute the difference in variance between the two alternatives and take the square root of this difference. The result thus achieved for each industry can be interpreted as the "room" we have for a model-based estimator; an RMSE (root mean squared error) for that estimator. When using stratified sampling we need a minimum of observations in each stratum. In our example, we get more than 400 "extra" observations for the full population already with at least 3 observations per stratum, and 100 of these extra observations are in size-class 0. When we delete this size-class, the number of "extra" observations decreases by 120. The lengths of the confidence intervals are roughly 90 % of the previous lengths. The room for the model-based estimator is considerable in most industries, several hundred percent in many cases. However, the room is smaller for some industries and, of course, overall. The figures here are a bit optimistic, since we have removed the extreme cases.

The effects of the second step – when we cut off also the next size-class with 1-4 employees – are similar to those of the first step in many respects, but not all. The sampling variance decreases considerably, and the confidence intervals are roughly 80 % of the previous lengths. There is more room for the model-based estimator than before in value, but much less in percent. The room is not worryingly small in most industries, but for a few industries a further check is motivated.

In the third step – up to the present cut-off threshold – the findings are like those in the second step but more pronounced. The differences between industries are considerable; the threshold at 10 employees is high for some but lower for others. The pressure for an accurate model varies accordingly. Some industries on the two-digit level have more than 10 % of the turnover below the threshold. On the other hand, these industries mostly have small shares of the overall total. For other industries, the share below the threshold is small, 1 % or less. We can compare the results here with those in Section 4. The three industries with NACE codes 18, 19, and 28 have high proportions of turnover also below the cut-off threshold. The two industries with NACE codes 21 and 34 are among the industries with low shares of small enterprises and much "room" for model-based estimation. For these two industries this is the case also if the threshold is further increased to 20 employees. However, in such a step the room is small for most industries.

## 7. CONCLUSIONS

We have gained some experience, and we find cut-off sampling to be a useful method – but a method that needs careful preparation before it is used. It is important to have an appropriate measure of size. We have found the variable "number of employees" too weak on its own. Especially enterprises without employees are a heterogeneous group, and a further grouping is needed. Other size measures should be investigated. The survey design should take the findings for size into account, whether cut-off sampling is used or not.

When a cut-off threshold is introduced it should be tailored to the survey and its focus. It is not appropriate to use one and the same threshold for all industries; we have seen important differences between industries on the two-digit level of the NACE code. For our case study, the present threshold can be increased in a few carefully selected industries, but it should be kept, or even decreased, for other industries.

It is not trivial to find an appropriate model for estimates for the population part that has been cut off. There are differences between size-classes. We have seen this for both the very small and the fairly small enterprises. However, some relationships between variables seem to be (fairly) stable over time, and such relationships can improve the model-based estimator. The model then uses both current information from the sample and earlier relationships from administrative data. Some effort is needed to build such models, which may have industry- and size-specific components, and there should be a regular supervision. The cost of these efforts is a part of the total costs to be balanced with accuracy.

## REFERENCES

Elvers, E. (1993) "A New Swedish Business Register Covering a Calendar Year and Examples of its Use for Estimation", *Proceedings of the International Conference on Establishment Surveys, American Statistical Association*, pp. 916-919.

Haan, J. De, E. Opperdoes, and C.M. Schut (1999). "Item Selection in the Consumer Price Index: Cut-off Versus Probability Sampling", *Survey methodology*, 25, pp. 31-41.

Särndal, C.-E., B. Swensson, and J. Wretman (1992) *Model Assisted Survey Sampling*, New York: Springer-Verlag.

Statistics Canada (2001) "Monthly Survey of Manufacturing (MSM)", *Statistical Data Documentation System*, Reference Number 2101, Statistics Canada.