

# A Theory of Quota Surveys

JEAN-CLAUDE DEVILLE<sup>1</sup>

## ABSTRACT

Simple or marginal quota surveys are analyzed using two methods: (1) behaviour modelling (super-population model) and prediction estimation, and (2) sample modelling (simple restricted random sampling) and estimation derived from the sample distribution. In both cases the limitations of the theory used to establish the variance formulas and estimates when measuring totals are described. An extension of the quota method (non-proportional quotas) is also briefly described and analyzed. In some cases, this may provide a very significant improvement in survey precision. The advantages of the quota method are compared with those of random sampling. The latter remains indispensable in the case of large scale surveys within the framework of Official Statistics.

**KEY WORDS:** Quota surveys; Super-population models; Restricted sampling; Regression estimation.

## 1. INTRODUCTION

Quota sampling is the method most frequently used in France by private polling institutions. It is easy to implement, inexpensive, and has many practical advantages. However, its disadvantages are also well known: likelihood of bias, no possibility of processing non-responses, and the need for external information in order to set the quotas. In the English literature (Cochran 1977; or Madow *et al.* 1983, for example) quotas have a very bad reputation due to the lack of a reliable theory on which statistical inference can be based. The only “defenders” of the method (Smith 1983, in particular) base their arguments on the principles of inference conditional upon sampling, where the sampling plan may generally be ignored.

This paper proposes a theory of quota surveys based on two types of modelling: population behaviour modelling (which is the approach of Smith or the ideas expressed in Gourieroux 1981), and modelling the method of sample collection, which may correspond to a more realistic idea.

In both cases, variance estimates are obtained by resorting to variations of regression estimators.

The first section of the paper describes the quota method and the results of the survey theory that can be subsequently useful. Parts 2 and 3 develop models for the behaviour of individuals in a population, or of those conducting the survey, which justify the method. The last section examines the problems raised, and attempts to demonstrate how the quota method can be used to add to the traditional probabilistic methods, rather than compete with them.

## 2. A BRIEF REVIEW OF THE QUOTA METHOD AND SURVEY THEORY

### 2.1 Cell Quotas; Quotas on the Margins of a Contingency Table – Some Practical Aspects of the Method

At the simplest level, the quota method resembles stratified sampling. The distribution in the population of a discrete characteristic  $h$  possessed by  $N_h$  individuals ( $h = 1$  to  $H$ ) is known.

<sup>1</sup> Jean-Claude Deville, Institut National de la Statistique et des Études Économiques, 18, Boulevard Adolphe Pinard, 75675, Paris Cedex 14, France.

The sample includes  $n_h$  individuals in category  $h$ ; however, the choice of these individuals is left up to the those conducting the survey. The sampling rate  $f_h = n_h/N_h$  may vary from category to category.

In practice, we prefer to control several criteria expressed as  $i, j, \dots, h$  ( $i = 1$  to  $I$ ,  $j = 1$  to  $J$ ,  $\dots$ ,  $h = 1$  to  $H$ ). Ideally, knowing the  $N_{ij\dots h}$  values of the multiple-entry contingency table allows the use of the previous method to define the number  $n_{ij\dots h}$  of members in the sample depending upon the  $f_{ij\dots h}$  rates. Except in very specific cases (few criteria having few modalities each) this method is unrealistic, because it leads to a search for individuals who are extremely difficult to find.

Thus, it is preferable to use **marginal quotas**, by calibrating the sample so that its distribution in accordance with the first criterion leads to a given  $n_{i+ \dots +}$  number of members, and the same is done for the other criteria. The only constraint on these marginal values is that they must be added to the overall sample size  $n$ . However, in practice, a single sampling rate  $f$  is adopted for each set of quotas:  $n_{i+ \dots +} = fN_{i+ \dots +}$ ,  $n_{+j \dots +} = fN_{+j \dots +}$  and  $n_{+ \dots +h} = fN_{+ \dots +h}$  with the obvious notations (+ in place of an index indicates the addition of all the modalities in the category represented by the index).

Beyond the obvious collection advantages, this technique is the one most often imposed by the external data on which the quotas are based. These are obtained, for example, from various sources, thus preventing any cross-correlations. Another situation arises when the quotas are established on the basis of a large survey (a labour survey, for example): each distribution is done in accordance with a criterion (age, socio-professional category, *etc.*) that may be considered to be reliable. On the other hand, the cross-correlations are affected by a large random error, and cannot be used to set the quotas.

In practice, the quota method is most often used to complement more traditional methods as the last sampling technique used in a multi-stage stratified survey on a geographic basis (region, size of the agglomerations). Each primary unit is assigned to a survey officer for whom quotas have been set. The survey officer also receives instructions to distribute his sample in order to make data collection as close to random as possible.

## 2.2 Traditional Survey Theory

We want to measure the total  $Y$  of a variable whose value  $Y_k$  for individual  $k$  is fixed, with no randomness. Only sample  $s$  is random, and the law of probability that governs  $s$  is known, since it is controlled by the statistician. Thus, we also know the possibility  $\pi_k$  that each individual will appear in  $s$ . Without any other information, the natural (unbiased) estimator to be used is the estimator based on inflated values:

$$\hat{Y} = \sum_{k \in s} Y_k / \pi_k = \sum_s d_k Y_k \quad \text{with} \quad d_k = 1 / \pi_k.$$

When the  $\pi_k$  are all equal to  $n/N$ , the sampling rate, we have:

$$\hat{Y} = N/n \sum_s Y_k = N\bar{y},$$

where  $\bar{y}$  is the mean of  $Y$  in the sample.

This estimator has a known variance, which is a quadratic form  $V(Y_U)$  on the vector of  $Y_k$  in the population:

$$\text{Var}(\hat{Y}) = V(Y_U) = \sum_k Y_k(d_k - 1) + \sum_{kl} Y_k Y_l d_k d_l (\pi_{kl} - \pi_k \pi_l), \quad (2.2.1)$$

where  $\pi_{kl}$  is the probability of simultaneously having  $k$  and  $l$  in  $s$ .

Similarly, the variance of  $\hat{Y}$  can be estimated by a quadratic form on vector  $Y_s$  of the  $Y_k$  in the sample:

$$\hat{V}(Y_s) = \sum_{kl \in s} \Delta_{kl} Y_k Y_l,$$

with

$$\begin{aligned} \Delta_{kl} &= (1 - \pi_k) / \pi_k^2 \quad \text{if } k = l \\ &= (\pi_{kl} - \pi_k \pi_l) / (\pi_{kl} \pi_k \pi_l) \quad \text{if } k \neq l. \end{aligned}$$

Depending upon the sampling plans, these expressions take the specific forms found in the manuals (Desabie 1965; Cochran 1977; Wolter 1985).

Any external information can improve the quality of the estimate. This is usually presented in the form of a vector  $X$  in which each of the  $p$  components is the total of a measurable variable in each of the possible samples. The estimate of  $Y$  can thus be improved by using regression estimation:

$$\hat{Y}_{\text{Reg}} = \hat{Y} + (X - \hat{X})' \hat{B},$$

where  $B$  is the vector of the coefficients of the regression of the  $Y_k$  on the  $X_k$  estimated by:

$$\hat{B} = \sum_s (d_k X_k X_k')^{-1} \sum_s d_k X_k Y_k.$$

When the constant is part of the regressors, or if it is a linear combination of the regressors and the sample has equal probabilities, the formula is simplified as follows:

$$\hat{Y}_{\text{Reg}} = X' \hat{B}.$$

The variance of  $\hat{Y}_{\text{Reg}}$  is simply expressed by introducing the residuals of the regression  $E_k = Y_k - X_k' B$  into the population. We know that we have:

$$\text{Var}(\hat{Y}_{\text{Reg}}) = V(E_U)$$

thus, we introduce in formula (2.2.1) vector  $E_U$  of residuals  $E_k$ . At the same time, we approximate an estimate of this variance by  $\hat{V}(e_s)$ , where  $e_s$  is the vector of  $e_k = Y_k - X_k' \hat{B}$ , the estimated residuals of the regression.

Under some sampling plans, these expressions assume particular forms. As a general rule,  $V$  and  $\hat{V}$  are the positive quadratic forms, and the  $E_k$  or  $e_k$  quantities smaller than the  $Y_k$ ; the regression estimator leads to substantial improvements over the inflated values.

A particularly important case that we will use later is one where  $X$  is a vector of the total accounting variables (values on the basis of which the quotas are constructed). Typically, the additional information is the vector of dimension  $I + (J - 1) + \dots + (H - 1)$  formed by the quantities:  $N_{i+\dots+}, N_{+j+\dots+}, N_{+\dots+h}$  for  $i = 1$  to  $I, j = 1$  to  $J - 1$ , and  $h = 1$  to  $H - 1$  (keeping only those variables that are linearly independent). Thus, the regressors

are the indicative variables of categories  $i$  ( $i = 1$  to  $I$ ),  $j$  ( $j = 1$  to  $J - 1$ ), and  $h = 1$  to  $(H - 1)$ . Since the constant is a linear combination of the regressors (it is the sum of the first  $I$  of them), the regression estimator takes the form:

$$\hat{Y}_{\text{Reg}} = \sum_i N_{i+} \hat{A}_i + \sum_j N_{+j} \hat{B}_j + \dots + \sum_h N_{+ \dots h} \hat{C}_h, \quad (2.2.2)$$

where  $\hat{A}_i$  (for example) indicates belonging to category  $i$ .

If we are only working with a single category, the regressors are orthogonal 2 by 2 and we have:

$$\hat{Y}_{\text{Reg}} = \sum_i N_i \hat{Y}_i$$

where  $\hat{Y}_i$  is the estimator of the mean of  $Y$  in category  $i$ . Thus,  $i$ .  $\hat{Y}_{\text{Reg}}$  is nothing but the post-stratified estimator.

### 2.3 Sampling Theories Based on Models

In this approach, we consider that the  $Y_k$  are random variables governed by a super-population model. This consists of parameters that we estimate on the basis of the sample. We can then calculate the probability, under the estimated model, of the non-observed values of  $Y$ , that is,  $\hat{Y}_k$ . The prediction estimator is the sum of the observed and predicted values and can be obtained as follows:

$$\hat{Y}_{\text{Pred}} = \sum_s Y_k + \sum_{U-s} \hat{Y}_k.$$

If, for example, in an equal probabilities survey, the model is a regression  $Y_k = X'_k \cdot \beta + \epsilon_k$ ,  $\epsilon_k$ , when the  $k$  values are independent, centred, and of equal variance, and when the constant appears on the regression (or when we have a linear combination of  $X_k$  that is constant), we have:  $\sum_s Y_k = \sum_s X'_k \beta$ ; and the prediction estimator and the regression estimator are the same.

We say that  $\hat{Y}$  is without bias under the model when, for all  $s$ ,  $\mathcal{E}(\hat{Y} - Y) = 0$  (conditionally upon the sample, the probability and variance under the model are expressed as  $\mathcal{E}$  and  $\mathcal{V}$ ). For the prediction estimator, we must only have, for all  $k$ , the natural condition  $\mathcal{E} \hat{Y}_k = \mathcal{E} Y_k$ , in order for this to be true. With the model, we can also evaluate the average quadratic deviation:  $\mathcal{E}(\hat{Y}_{\text{Pred}} - Y)^2$ , since we know that the two terms  $\hat{Y}_{\text{Pred}}$  and  $Y$  are random, and that  $\hat{Y}_{\text{Pred}}$  depends upon sample  $s$ . The above-mentioned probability is thus conditional upon sample  $s$ . This follows a certain probability law already discussed in the previous paragraph. The precision of this estimator can be measured by calculating:

$$\mathcal{V}(\hat{Y}_{\text{Pred}}) = E\mathcal{E}(\hat{Y}_{\text{Pred}} - Y)^2.$$

If the law of  $s$  is such that the  $Y_k$  are independent (the so-called non-informative sampling), then this quantity equals:

$$\mathcal{E}(E(\hat{Y}_{\text{Pred}} - Y)^2),$$

where the internal probability is conditional upon  $Y_k$ . If  $\hat{Y}_{\text{Pred}}$  is equal to  $\hat{Y}_{\text{Reg}}$ , and we have a condition of independence, we will have:

$$\mathcal{V}(\hat{Y}_{\text{Pred}}) = \mathcal{E}(\text{Var}(\hat{Y}_{\text{Reg}})).$$

## 2.4 Comments on the Two Approaches Applied to the Quota Method

a) In both cases, the process of estimation will be effective if the variable of interest is well explained by category indicators on which the quotas are roughly based, because the regression adjustment residuals will be small.

b) In a quota survey the “sampling plan” is not known by the statistician. Thus, he cannot make inferences without using a model. The latter may be a population behaviour model (“model” approach) that requires him to assume certain responsibilities regarding the nature of what he observes. This approach will be developed in the second part of this paper. This may also consist of modelling the sampling plan; which means taking responsibility for the operation of the collection process. This approach will be developed in the third section of this paper.

In all cases, the modelling speculation must be mobilized in order to validate a kind of inference. The question is to know whether it is easier and more plausible to model the behaviour of the individuals surveyed, or to model the sample collection process (including the contacts between interviewer and interviewee).

c) In this respect, the hypothesis made in section 2.3 regarding the independence between randomness in the population and randomness in the collection process is **crucial**. If sampling is controlled by the statisticians, this guarantee can be ensured, except for the effect of non-responses. In the case of the quota method, there are no guarantees. Let us assume, for example, that we want to measure incomes  $Y_k$ , the probability  $\pi_k$  of finding  $k$  in the sample may be very low if  $Y_k$  is large. In other words, the fact of belonging to the sample (which is 1 if  $k$  is in  $s$ , and 0 otherwise) and the residual of the super-population model  $\epsilon_k$  are negatively correlated. This example illustrates well the main danger of the quota method, which the following theory does not take into account.

## 3. QUOTA THEORY WITH A SUPER-POPULATION MODEL

### 3.1 Cell Quotas

There is a single cell category  $i = 1$  to  $I$  for the known values  $N_i$ . The model that can be imagined is as follows:

$$Y_k = m_i + \epsilon_k, \quad (3.1.1)$$

$\epsilon_k$  centred independently of variance  $\sigma_i^2$  where  $i$  is the cell to which  $k$  belongs.

The Gauss-Markov estimators of  $m_i$  are the means observed in the various  $\bar{y}_i$  cells. Thus, the prediction estimator is:

$$\hat{Y}_{\text{Pred}} = \sum_i (N_i - n_i) \bar{y}_i + \sum_i n_i \bar{y}_i = \sum_i N_i \bar{y}_i. \quad (3.1.2)$$

This has the form of the post-stratified estimator. Moreover:

$$\text{Var}(\hat{Y}_{\text{Pred}} - Y)^2 = \sum_i \sigma_i^2 N_i (N_i - n_i) / n_i. \quad (3.1.3)$$

This quantity does not depend upon sample  $s$ , as the latter always includes (with a probability of 1 !)  $n_i$  individuals in cell  $i$ .

$E\mathcal{E}(\hat{Y}_{\text{Pred}} - Y)^2$  can be estimated by replacing  $\sigma_i^2$  by its usual estimator  $s_i^2 = (n_i - 1)^{-1} \sum_{k \in s_i} (Y_k - \bar{y}_i)^2$  with  $s_i$  being part of  $s$  in cell  $i$ .

These results are from Gourieroux (1981) and represent, to a certain extent, a justification of the simple quota method.

### 3.2 Marginal Quotas – “Representative” Case

In this and the following paragraphs, we will restrict ourselves to the case of quotas overlapping 2 criteria  $i$  and  $j$ . The generalization with more than 2 criteria does not pose any particular problems, but leads to very complex notations that we prefer to avoid (see Appendix).

Thus, the situation is as follows: the values  $N_{i+}$  and  $N_{+j}$  of the two universe breakdowns are known. The sampling only allows samples of fixed size  $n = fN$  including  $n_{i+} = fN_{i+}$  individuals for each  $i$ , and  $n_{+j} = fN_{+j}$  individuals for each  $j$ .

We postulate an analysis of variance model in the population, formulated as follows:

If  $k$  belongs to cell  $(i, j)$ :

$$Y_k = \alpha_i + \beta_j + \epsilon_k. \tag{3.2.1}$$

The  $\epsilon_k$  are centred, independent, and we have  $\text{Var } \epsilon_k = \sigma_i^2 + \gamma_j^2$ .

For reasons of identification of the model, we postulate that  $\beta_J = 0$ .

This is equivalent to postulating that  $Y_k = (\alpha_i + u_{ik}) + (\beta_j + v_{jk})$  where  $u_{ik}$  and  $v_{jk}$  are independent, and their respective variances are  $\sigma_i^2$  and  $\tau_j^2$ .

We estimate  $\alpha_i$  and  $\beta_j$  using the ordinary least squares (OLS) method, because we ignore the values of the variance elements; the  $\hat{\alpha}_i$  and  $\hat{\beta}_j$  are solutions of the system:

$$\begin{aligned} \sum_j n_{ij} \bar{y}_{ij} &= n_{i+} \hat{\alpha}_i + \sum_j n_{ij} \hat{\beta}_j \quad (i = 1 \text{ to } I) \\ \sum_i n_{ij} \bar{y}_{ij} &= n_{+j} \hat{\beta}_j + \sum_i n_{ij} \hat{\alpha}_i \quad (j = 1 \text{ to } J - 1), \end{aligned} \tag{3.2.2}$$

with  $\bar{y}_{ij}$  the mean of the  $Y_k$  over the  $s_{ij}$  part of the sample in cell  $(i, j)$ . Thus, the prediction estimator can be written as follows:

$$\hat{Y}_{\text{Pred}} = \sum_{ij} (N_{ij} - n_{ij}) (\hat{\alpha}_i + \hat{\beta}_j) + \sum_{ij} n_{ij} \bar{y}_{ij}.$$

**Result 1:** Under model (3.2.1), the prediction estimator using the OLS is  $N\bar{y}$ . We check that it is unbiased for the model; that is, that  $\mathcal{E}(N\bar{y} - Y) = 0$ .

**Proof:** Immediately from (3.2.2), and because of the fact that the quotas are proportional to the numbers in the population.

**Result 2:** We have:

$$\mathcal{E}(N\bar{y} - Y)^2 = (N^2/n)(1 - f)n^{-1} \left( \sum_i n_{i+} \sigma_i^2 + \sum_j n_{+j} \tau_j^2 \right).$$

This quantity does not depend upon the sample (as it depends only upon the quotas). Thus, to a certain extent, this is a justification for the marginal quotas method.

**Proof:** With  $m_k = \varepsilon Y_k$ , using the unbiased character of the estimator we have:

$$\begin{aligned} \varepsilon(N\bar{y} - Y)^2 &= \varepsilon\left((N/n) \sum_s (Y_k - m_k) - \sum_U (Y_l - m_l)\right)^2 \\ &= \varepsilon\left((N/n) \sum_s \epsilon_k - \sum_U \epsilon_l\right)^2 \\ &= (N/n)^2 \sum_{ij} n_{ij}(\sigma_i^2 + \tau_j^2) - 2(N/n) \sum_{ij} n_{ij}(\sigma_i^2 + \tau_j^2) + \sum_{ij} N_{ij}(\sigma_i^2 + \tau_j^2). \end{aligned}$$

But

$$\begin{aligned} \sum_{ij} N_{ij}(\sigma_i^2 + \tau_j^2) &= \sum_i N_{i+} \sigma_i^2 + \sum_j N_{+j} \tau_j^2 \\ &= (N/n) \left( \sum_i n_{i+} \sigma_i^2 + \sum_j n_{+j} \tau_j^2 \right) \end{aligned}$$

from which:

$$\begin{aligned} \varepsilon(N\bar{y} - Y)^2 &= (N^2/n)(1 - f)n^{-1} \left( \sum_{ij} n_{ij}(\sigma_i^2 + \tau_j^2) \right) \\ &= (N^2/n)(1 - f) \left( \sum_i p_{i+} \sigma_i^2 + \sum_j p_{+j} \tau_j^2 \right) \end{aligned}$$

$$\text{with } p_{i+} = N_{i+}/N \text{ and } p_{+j} = N_{+j}/N.$$

The estimate of the precision of  $E(N\bar{y} - Y)^2$  is derived from this. In fact, with this model,  $s_{ij}^2$  has a probability of  $\sigma_i^2 + \tau_j^2$ . Thus, an unbiased estimator of the precision is obtained by

$$(N/n)^2 (1 - f) \sum_{ij} n_{ij} s_{ij}^2$$

if all the  $n_{ij}$  are equal to or greater than 2.

This estimator is formally identical to the one that we would use in a complete post-stratification on cells  $(i, j)$ . We can also use  $(N/n)^2 (1 - f) \sum_s e_k^2$ , where  $e_k$  are the estimated residuals of the model.

### 3.3 What Happens if the Model is False?

**3.3.1** An initial way of looking at the question is to put model (3.2.1) into the general model where the mean of  $Y_k$  depends upon the pair  $(i, j)$ . This can be written as follows:

$$Y_k = \alpha_i + \beta_j + \gamma_{ij} + \epsilon_k, \tag{3.3.1.1}$$

with the usual hypotheses for  $\epsilon_k$  and the terms of interaction  $\gamma_{ij}$  that verify the constraints of identifiability:

$$\sum_j N_{ij} \gamma_{ij} = 0 \text{ and } \sum_i N_{ij} \gamma_{ij} = 0. \tag{3.3.1.2}$$

Thus we have:

$$\varepsilon(N\bar{y} - Y) = \sum_{ij} (Nn_{ij}/n - N_{ij})\gamma_{ij}, \quad (3.3.1.3)$$

such that the estimator is biased for the model except when  $n_{ij} = fN_{ij}$ , which has no reason to exist.

This means that the terms of sum (3.3.1.3) may well compensate for each other, since their signs are *a priori* undetermined.

On the other hand, if “good” sampling precautions are taken,  $Nn_{ij}/n - N_{ij}$  should usually be close to 0.

It is clear, in any case, that the more suitable the additive model is (small  $\gamma_{ij}$ ), and the more the sampling plan approaches randomness, the more likely it is that bias will be reduced.

**3.3.2** Another way to view the misrepresentation of the model, which has already been described, is to no longer admit that there is independence between the randomness of the sample and the randomness of the additive model. This means that distinct models should be developed for the  $(Y_k, k \in S)$  and  $(Y_l, l \notin S)$  vectors. This approach has often been used in the econometric literature, to which the reader is referred. It is clear that risk-taking in regards to the data becomes enormous, and is often incompatible with objective work on the part of the statistician.

### 3.4 Marginal Quotas with Unequal Rates

In the case of cell quotas, we can arbitrarily set quotas for each cell. Until now, in the case of marginal quotas, we have only examined the case where the quotas were proportional to the size of the population.

In many cases however, we may be tempted to over-represent certain categories. If, for example, we want to study household assets, we may want to set the largest quotas for older households (quotas by age group), on the one hand; and for those where the head is self-employed (quotas by social categories), on the other.

Thus, we formally force the sample to fall within a given size  $n_{i+}$  and  $n_{+j}$  (however, the sum of  $n_{i+}$  is always equal to the sum of  $n_{+j}$ ).

In this case, always using the OLS as an estimation technique, we can easily find that the total prediction estimator is:

$$\hat{Y}_{\text{Pred}} = \sum_i N_{i+} \hat{\alpha}_i + \sum_j N_{+j} \hat{\beta}_j, \quad (3.4.1)$$

$\hat{\alpha}_i$  and  $\hat{\beta}_j$  always verify estimating equations (3.2.2). It is easy to see that this estimator may be expressed as follows:

$$\hat{Y}_{\text{Pred}} = \sum_{ij} (w_i^{(1)} + w_j^{(2)}) n_{ij} \bar{y}_{ij} = \sum_{ij} \hat{N}_{ij} \bar{y}_{ij}.$$

Thus, the quantities  $(w_i^{(1)} + w_j^{(2)}) n_{ij}$  seem to be estimates of the size of cells  $(i, j)$ , an idea that will be largely exploited in the following sections.

On the other hand, the variance of this estimator under the model depends upon all the  $n_{ij}$ , and this can be demonstrated by a rather cumbersome calculation. The justification of the quota method described above no longer works.



#### 4. MODELS FOR THE SAMPLING PLAN

##### 4.1 A Model Sampling Plan

The idea is one of a simple random sampling constrained by the quotas imposed. The selection algorithm, while totally unrealistic, consists of drawing a series of simple random samples until we find one that verifies the quotas. Thus, each sample that verifies the quotas has the same positive probability of being drawn, the samples that do not verify the quotas have a zero probability of being drawn.

The purpose is to model the fact that the person conducting the survey will correctly follow the dispersion constraints on the survey units assigned to him.

##### 4.2 Cell Quotas

This sampling model is based on an *a priori* stratification. Its practical advantage is that it does not require a sampling frame where the stratification variables are present. It is implemented rigorously in certain cases, for example, in a telephone survey based on a non-informative random list of telephone numbers, and when surveys are carried out only until the quotas are met.

The formulas that provide the estimators, the variances, and the precision estimates are those given in all the manuals. They have a certain similarity with those described in section 3.1 (see Gouriboux 1981).

##### 4.3 The Case of Marginal Quotas: General Estimators

The sampling model is that of simple random sampling constrained by marginal quotas. SRS provides samples with  $n_{ij}$  members in the various cells that can be taken as a random vector (in whole values) in  $R^{IJ}$ . The quota constraint means that we are limited to a random vector as follows:

$$\sum_j n_{ij} = n_{i+} \quad (i = 1 \text{ to } I) \quad \text{and} \quad \sum_i n_{ij} = n_{+j} \quad (j = 1 \text{ to } J - 1),$$

that is, one that varies within a sub-space of size  $IJ - I - J + 1$ . We place ourselves in the case where the overall sampling rate is negligible, and the law of the  $n_{ij}$  can be compared to a multinomial law ( $n, p_{ij} = N_{ij}/N$ ).

Conditional upon  $n_{ij}$ , the  $\bar{y}_{ij}$  estimate the  $\bar{Y}_{ij}$  without bias. The idea is now to construct an estimator of the total of  $Y$  by weighting the  $\bar{y}_{ij}$  by the estimators of  $N_{ij}$ , that is, the  $p_{ij}$ . If we choose to maximize the probability, this is proportional to:

$$\prod_{ij} p_{ij}^{n_{ij}}. \tag{4.3.1}$$

Thus, we maximize

$$\sum_{ij} n_{ij} \text{Log} p_{ij} \tag{4.3.2}$$

under the following constraints

$$\sum_j p_{ij} = p_{i+} \quad (i = 1 \text{ to } I) \quad \text{and} \quad \sum_i p_{ij} = p_{+j} \quad (j = 1 \text{ to } J - 1) \tag{4.3.3}$$

which leads to solving the system for  $a_i, b_j$  ( $p_{i+} = N_{i+}/N, p_{+j} = N_{+j}/N$  are known):

$$\sum_j \hat{p}_{ij}^{\circ} (a_i + b_j)^{-1} = p_{i+} \quad (i = 1 \text{ to } I) \quad (4.3.4)$$

$$\sum_i \hat{p}_{ij}^{\circ} (a_i + b_j)^{-1} = p_{+j} \quad (j = 1 \text{ to } J - 1; b_J = 0),$$

with  $\hat{p}_{ij}^{\circ} = n_{ij}/n$  frequency in the sample.

The estimators of  $p_{ij}$  are thus  $\hat{p}_{ij}^{\circ} (a_i + b_j)^{-1}$  and the estimator we are looking for can be written as follows:

$$\hat{Y}_Q = (N/n) \sum_{ij} n_{ij} (a_i + b_j)^{-1} \bar{y}_{ij} = (N/n) \sum_s w_k Y_k, \quad (4.3.5)$$

where  $w_k = (a_i + b_j)^{-1}$  is the weight added to  $Y_k$  in the case when  $k$  appears in cell  $(i, j)$ . This estimator is asymptotically without bias under the SRS model in  $U$ , as are the maximum probability estimators. The quotas do not play an explicit role in (3.3.4), but they affect the values of  $a_i$  and  $b_j$ .

In the normal case when the marginal quotas are ‘‘proportional’’, with a fixed sampling fraction  $f$ , the solution of equations (4.3.4) is evident:  $a_i = 1$  for any  $i$ , and  $b_j = 0$  for any  $j$ . The estimator of the total is  $N\bar{y}$ , as could be expected, and has the same expression as the equal-probability probabilistic sampling.

**Comment:** The use of maximum probability to estimate the proportions is rather arbitrary. A chi-square criterion (minimize  $\sum_{ij} (p_{ij} - \hat{p}_{ij}^{\circ})^2 / \hat{p}_{ij}^{\circ}$ ) would make the (4.3.4) system linear.

#### 4.4 Variance of the Estimator and its Estimate

**4.4.1** To establish a variance formula we will use the parametrization of variable  $Y$  used by J.C. Deville and C.E. Särndal (1990), which we will express in the form of a:

**Lemma:** For any variable  $Y = (Y_k; k \in U)$ , we can choose an uniquely defined parametrization

$$Y_k = \bar{Y}_{ij} + R_k \quad \text{if } k \text{ is in cell } (i, j) \quad (k \in U_{ij}) \quad \text{with} \quad \sum_{k \in U_{ij}} R_k = 0,$$

$$\bar{Y}_{ij} = A_i + B_j + E_{ij} \quad \text{with} \quad B_J = 0$$

$$\sum_j N_{ij} E_{ij} = 0 \quad i = 1 \text{ to } I$$

$$\sum_i N_{ij} E_{ij} = 0 \quad j = 1 \text{ to } J - I.$$

In fact,  $A_i$  and  $B_j$  are numbers that minimize the quantity  $\sum_U (Y_k - A_i - B_j)^2$  where, in an equivalent manner  $\sum_{ij} N_{ij} (\bar{Y}_{ij} - A_i - B_j)^2$ .

Thus, we can write:

$$\hat{Y}_Q = (N/n) \sum_{ij} n_{ij} (a_i + b_j)^{-1} (A_i + B_j + E_{ij} + \bar{R}_{ij}) \quad \text{where} \quad \bar{R}_{ij} = \sum_{s_{ij}} R_k/n_{ij}.$$

Taking into account equation 4.3.4 and the lemma:

$$\hat{Y}_Q - Y = \sum_{ij} \hat{N}_{ij} (E_{ij} + \bar{R}_{ij}) \quad \text{with} \quad \hat{N}_{ij} = (N/n) n_{ij} (a_i + b_j)^{-1}, \quad (4.4.1)$$

which is the basic expression for the calculation of the variance.

Conditional upon  $n_{ij}$ , the  $\hat{N}_{ij}$  are constant, and sub-samples  $s_{ij}$  are independent simple random samplings. Thus we have:

$$\text{Cond bias}(\hat{Y}_Q) = \sum_{ij} \hat{N}_{ij} E_{ij} = N \sum_{ij} \hat{p}_{ij} E_{ij}$$

$$\text{Cond Var}(\hat{Y}_Q) = \sum_{ij} \hat{N}_{ij}^2 V_{ij}/n_{ij} \quad \text{where} \quad V_{ij} = (1/N_{ij}) \sum_{U_{ij}} R_k^2.$$

Thus (demonstration in the Appendix) we have:

**Result 1:**

$$\text{Var} \left( \sum_{ij} \hat{p}_{ij} E_{ij} \right) = 1/n \sum_{ij} p_{ij} E_{ij}^2.$$

Furthermore, the probability of  $\hat{p}_{ij}^{\circ}(a_i + b_j)^{-1}$  is (in terms close to  $1/n$ )  $p_{ij}(a_i^{\circ} + b_j^{\circ})^{-1}$  where  $a_i^{\circ}$  and  $b_j^{\circ}$  are the solutions to equations (4.3.4), in which  $\hat{p}_{ij}^{\circ}$  are replaced by the exact  $p_{ij}$ .

This leads to:

**Result 2:** The variance of the quota estimator  $\hat{Y}_Q$  is given by:

$$\text{Var}(\hat{Y}_Q) = (N^2/n) \sum_{ij} p_{ij} (E_{ij}^2 + (a_i^{\circ} + b_j^{\circ})^{-1} V_{ij}).$$

If the quotas are proportional to the size of the population, we will have:

$$\text{Var}(\hat{Y}_Q) = (N^2/n) \sum_{ij} p_{ij} (E_{ij}^2 + V_{ij}).$$

#### 4.4.2 Estimating the Variance

The conditional variance of  $\hat{Y}_Q$  can be estimated by:

$$\sum_{ij} \hat{N}_{ij}^2 s_{ij}^2/n_{ij} = (N^2/n) \sum_{ij} \hat{p}_{ij} (a_i + b_j)^{-1} s_{ij}^2,$$

where  $s_{ij}^2$  is the usual unbiased estimator of  $V_{ij}$ . The probability of the square of the conditional bias is  $(N^2/n) \sum_{ij} p_{ij} E_{ij}^2$  and is estimated by  $(N^2/n) \sum_{ij} \hat{p}_{ij} \hat{E}_{ij}^2$  where  $\hat{E}_{ij} = \bar{y}_{ij} - \hat{A}_i - \hat{B}_j$  and  $\hat{A}_i$  and  $\hat{B}_j$  are the solutions of:

$$\sum_j \hat{p}_{ij} (\hat{A}_i + \hat{B}_j) = \sum_j \hat{p}_{ij} \bar{y}_{ij} \quad (i = 1 \text{ to } I), \quad (4.4.2)$$

$$\sum_j \hat{p}_{ij} (\hat{A}_i + \hat{B}_j) = \sum_j \hat{p}_{ij} \bar{y}_{ij} \quad (j = 1 \text{ to } J - I) \quad \text{with } B_J = 0.$$

In other words, the estimate of  $E_{ij}$  is obtained by fitting to the data an additive ANOVA model without interaction, the fitness criterion being that of least squares weighted by  $(a_i + b_j)^{-1}$ .

Thus, the variance estimator is:

$$\widehat{\text{Var}}(\hat{Y}_Q) = (N^2/n) \sum_{ij} \hat{p}_{ij} (\hat{E}_{ij}^2 + (a_i + b_j)^{-1} s_{ij}^2). \quad (4.4.3)$$

When the quotas are proportional to the population numbers, this expression can be simplified as follows:

$$(N^2/n) \sum_{ij} n_{ij} (\hat{E}_{ij}^2 + s_{ij}^2)/n. \quad (4.4.4)$$

If the  $n_{ij}$  are all sufficiently large that  $n_{ij}/(n_{ij} - 1) = 1$ , the sum of the formula is the sum of the squares of the residuals estimated in the OLS adjustment of the  $Y_k = A_i + B_j +$  residual model. Thus, the estimation procedure is simple:

- use the OLS to fit the additive model to the individual data
- create the variable  $e_k$  of the estimated residuals
- $\widehat{\text{Var}}(\hat{Y}_Q) = (N^2/n) \cdot (1/n) \sum_s e_k^2$ .

This formula is precisely that proposed in paragraph 2, and based on the super-population model. A rather neat situation!

#### 4.4.3 Discussion of the Results

The variance breaks down into two parts: one that can be seen as the probability of the square of the conditional bias; and one as the probability of the conditional variance.

The first term does not depend upon the quotas imposed on the sample, but only upon the quality of the fit of an additive model to the variable of interest. This part of the variance is diminished by choosing quota criteria that can best explain what we want to measure.

The second term, on the other hand, depends upon the remaining variability  $(N_{ij}^2 V_{ij}/n_{ij})$  and the number of observations collected in each cell. Since the size of the sample is fixed, we must attempt to make the  $n_{ij}$  as close as possible to Neyman's distribution:  $n_{ij} \propto N_{ij} V_{ij}^{1/2}$ . This may be achieved approximately by overloading quotas  $n_{i+}$  and  $n_{+j}$ , which correspond to large values of  $V_{ij}$ . Thus, in some cases, it is possible to improve the precision of a quota survey considerably.

## 4.5 Combination of the Quota Method and Stratified or Multi-Stage Samplings

### 4.5.1 The Case of Stratified Sampling with a Quota in Each Stratum

If the size of the criteria used to set the quotas are known in each stratum, the method described above makes it possible to construct an unbiased estimator, under the hypothesis that sampling functions like an SRS constraint in **each stratum**. If the allocation of quotas is proportional to the size of each stratum, the estimator is the natural estimator of the stratified sampling. If “national” quotas are used with each stratum, a correction should be made by reweighting.

On the other hand, if the size of the quota variables is unknown at the stratum level, it is not possible to correct the estimators to eliminate “structure effects” related to the stratification. Since, furthermore, the purpose of stratification is to construct dissimilar sub-populations, the corrections required will generally be quite large. Thus, the quota method is not recommended (except when the validity of the additive model is quite clear, *cf* part 3).

### 4.5.2 The Case of Two-Stage Sampling

Let us assume a two-stage sampling (inside a stratum where the sizes of the quota variables are known). If the sizes of the quota variables are known at the level of each primary unit, there are no problems. The theory in section 4.4 makes it possible to obtain an estimator of the total  $Y$  in each primary unit, as well as to calculate its variance, and an estimator of the latter. These quantities can then be used to obtain an estimator of  $Y$ , as well as an estimator of precision (*cf* Rao 1975). If the sizes of the quota criteria are not known at the level of the primary units, but only at the stratum level, we again have a problem that is impossible to correct. However, there is generally little harm if the PU are relatively similar: the structure of each PU is close to that of the stratum as a whole, and the corrections to be made for each PU are close to those that must be made at the stratum level.

### 4.5.3 In Conclusion

In conclusion, in the case complex multi-stage stratified sampling, the quota method may be used as the final sampling method if the stratification was carried out effectively by regrouping the similar primary units together, and if quotas derived from the data relative to each stratum are used with each PU.

To the extent that the hypothesis of simple random sampling constrained in each PU may appear to be quite satisfactory, the quota method is justified independently of any super-population model.

## 5. CONCLUSIONS AND PROBLEMS

### 5.1 How Should Non-response Be Taken into Account?

As we have already shown, this is the most important limitation in our theory. As far as sampling using the quota method is concerned, we do not have, in principle, any information on members of the population who refuse to respond to the survey, and we find ourselves lacking individual information on the subject of non-respondents. However, the situation is not as desperate as one might think. Let us illustrate this using a very simplified example.

We have carried out a simple quota survey using a sample of  $n_i$  individuals in category  $i$  with a population  $N_i$ . An acceptable model of non-response postulates a response probability of  $r_c$  if an individual belongs to category  $c$  with a population  $N_c$ . The (unknown) population

of the intersection between quota category  $i$  and class  $c$  of the non-response model is expressed as  $N_i^c$ . The population likely to respond in category  $i$  is thus  $N_{ri} = \sum_c N_i^c r_c$ . By setting a quota  $n_i$  in this category, within the framework of model (4.1), we obtain a probability of inclusion in the sample of  $w_i^{-1} = n_i/N_{ri}$ . In the sample, we collect  $n_i^c$  individuals belonging to the intersection  $(i,c)$  between the two categories. This quantity is random, and its probability is  $N_i^c r_c w_i^{-1}$ . If we attempt to estimate  $N_i^c$ , we will solve the estimating equations derived from the following relations:

$$N_i^c = n_i^c w_i r_c^{-1},$$

$$\sum_c N_i^c = N_i,$$

$$\sum_i N_i^c = N^c.$$

Thus, ranking ratio technique makes it possible to obtain estimates of  $f_c$  and  $\hat{w}_i$ , and to derive estimators  $\hat{N}_i^c = n_i^c \hat{w}_i \hat{r}_c^{-1}$  from the sizes of the intersection  $(i,c)$ . We can also obtain an estimator of the total of  $Y$ :

$$\hat{Y}_{NR} = \sum_{ic} N_i^c \bar{y}_i^c = \sum_{ic} r_c^{-1} w_i n_i^c \bar{y}_i^c,$$

where  $\bar{y}_i^c$  is the mean of the  $Y_k$  values in the sample in category  $(i,c)$ . Thus, estimation techniques based on fitting should allow for the honourable processing of non-responses in quota surveys.

## 5.2 Some Points of Comparison with Probabilistic Surveys

Regardless of how we try to understand it, the quota method demands the formulation of a hypothetical model to fit the data. On the other hand, a probabilistic survey does not, in principle, depend upon any model. In practice, sampling for a probabilistic survey is a model to which the reality of data collection attempts to conform. In fact, we are well aware that, in any probabilistic survey, some compromises of detail must be made with the model (necessary exclusion of certain units, replacement of others after selection but before data collection, *etc*). However, we can say that statistical biases are always much lower in probabilistic selection than when using the quota method. On the other hand, quotas make it possible to use, in the sampling stage, additional information that cannot be mobilized in a probabilistic selection process. As a result, the variance of a quota sampling is similar to that of a regression estimation, and is thus generally smaller than that resulting from a probabilistic survey associated with its estimate of standard inflated values. The choice is between bias due to the model associated with low variance, against lack of bias. Two types of conclusions can be drawn from this approach:

**5.2.1 Precision depends mostly upon the size of the sample.** On the average, in the case of small samples, probabilistic sampling will produce the worst results; and the bias of a quota survey will be more tolerable than the lack of precision of a probabilistic survey. For large samples, on the other hand, the quota method will have a clear bias that is obviously incompatible with the confidence interval without bias of a probabilistic survey.

Where should the boundary between the two methods be set? It is hard for the theory to be specific. On the other hand, experience in the French institutes may lead to a solution to this question: most national quota surveys are carried out on samples of 1,000 to 2,000 individuals. On the other hand, no national probabilistic survey mobilizes less than 5,000 units. It would seem fair to say that a size of 2,500 to 3,000 surveys is a practical boundary between the two types of surveys.

### **5.2.2 Official Statistics or Marketing**

In a survey, the use of any speculative model represents methodological risk-taking. This may be perfectly reasonable if the users are aware of it, and if they have ratified the speculations leading to the specification of the model. This is typically what happens, at least implicitly, in marketing surveys: an organization, company, administration, or association requests a sampling survey from a polling company. A contract marks the agreement between the two parties respecting the implementation of the survey, its price, the result delivery schedule, and **the methodology used**. In this methodology, models are used to formalize the sampling or behaviour of the population. Thus, from this point of view, the use of the quota method may be quite proper.

Official statisticians, on the other hand, are responsible for generating data that can be used by the entire society; and that can be used, in particular, in the arbitration of disputes between various groups, parties, and social classes. The use of statistical models, particularly econometric models that describe the behaviour of economic agents, may turn out to be very dangerous, partial, or affected by a questionable or disputed economic theory. Official statistics should not tolerate any uncontrollable bias in its products. It should carry out sample surveys using probabilistic methods.

There is no real opposition between quota survey techniques and those using controlled randomness, quite the opposite – they are complementary. As a proof of this, the statistics that are used to construct the quotas are themselves very often derived from large surveys carried out by the National Statistics Services. However, quota survey technicians find it hard to admit that these data are obtained using methods other than traditional, confirmed, and well-founded probabilistic techniques.

### **ACKNOWLEDGEMENTS**

I would like to express my sincere thanks to the referee and editor for their help in improving the quality of this paper.

## APPENDIX

### Demonstration of the Results of Section 4.4

#### 1. Notation and Results

In order to deal with the question in a general way, we will require certain convenient notations. We have  $Q$  qualitative variables whose modalities are indicated by using indices from 1 to  $I_q$  when  $q = 1$  to  $Q$ . A "cell" is denoted as  $c$ ; that is, a series of  $Q$  indices where the  $q^{\text{th}}$  could have a value of 1 to  $I_q$ ; and  $q_c$  is the value of the  $q^{\text{th}}$  index ( $q^{\text{th}}$  projection of  $c$ ); in a finite population  $U$ , of size  $N$ ,  $U_c$  is the population of individuals in cell  $c$ , when the size of the cell is  $N_c$ . The quantity  $N_i^{+q} = \sum_{q_c} i N_c$  is the total of the  $Q$ -dimensional contingency table where the cells are represented by  $c$  for the  $i^{\text{th}}$  modality of the  $q^{\text{th}}$  variable. If we postulate that

$$\bar{Y}_c = \frac{1}{N_c} \sum_{k \in U_c} Y_k.$$

We will obtain the following results:

**Result 1:** Variable  $Y_k (k \in U)$  may be parametrized by the following numbers:  $A_{q_c}^q$ ,  $E_c$  and  $R_k$  by:

$$\bar{Y}_k = \bar{Y}_c + R_k \quad \text{if } k \in U_c. \quad \text{We have } \sum_{U_c} R_k = 0 \quad \text{for any } c.$$

$$\bar{Y}_c = \sum_{q=1}^Q A_{q_c}^q + E_c \quad \text{with } A_{I_q}^q = 0 \quad \text{for } q = 2 \text{ to } Q \quad \text{and}$$

$$\sum_{q_c=i} N_c E_c = 0 \quad \text{for } q = i \text{ to } Q \quad \text{and } i = 1 \text{ to } I_q.$$

These numbers are obtained from the minimization of:

$$\sum_U \left( Y_k - \sum_{q=1}^Q A_{q_c(k)}^q \right)^2 = \sum_c N_c \left( \bar{Y}_c - \sum_{q=1}^Q A_{q_c}^q \right)^2.$$

Let us assume that we have a sample  $s$ . We will use  $n$  to denote all quantities in the sample that are similar to whatever we have already indicated in the population.

We assume that  $s$  was obtained on the basis of simple random sampling (with or without replacement) in accordance with an equal probability scheme constrained by the totals  $n_i^{+q}$  ( $q = 1$  to  $Q$ ,  $i = 1$  to  $I_q$ ), the quotas.

The purpose of this appendix is to demonstrate the following result:

**Result 2:** The variance of  $\sum_c \hat{N}_c E_c$  is approximately equal to  $1/n \sum_c N_c E_c^2$  when  $n$ , and  $N/n$  become arbitrarily large.

The following section will provide a more precise formulation for this result.

#### 2. Sampling Plan and Asymptotic Reduction

Let us consider the following two sampling models SR and AR:

**SR:** Bernoulli Sampling. Each of the units of  $N$  belong to  $s$  with a probability  $f$ , and the  $N$  drawings are independent.



AR: Each unit is drawn a number  $v_k$  of times;  $v_k$  follows Poisson's law with  $f$  parameters. The  $v_k$  are independent variables.

A simple random survey without replacement (SRSWOR) of fixed size  $n$  is an SR sampling if the total size of the sample is  $n$ .

A simple random survey with replacement (SRSWR) of fixed size  $n$  is an AR sampling when we have  $n$  observations; that is, when  $\sum_k v_k = n$ .

In the case of SR sampling, the law of the vector  $n_c$  is obtained as follows:

$$\Pr(\{n_c\}) = \prod_c \binom{N_c}{n_c} f^{n_c} (1 - f)^{N_c - n_c}.$$

In the case of AR sampling, we have:

$$\Pr(\{n_c\}) = \prod_c \frac{(N_c f)^{n_c}}{n_c!} \exp(-fN_c).$$

In both cases the variables  $n_c$  are independent.

In the case of SR sampling constrained by  $\sum n_c = n$ , the law of the  $n_c$  is hypergeometric:

$$\Pr(\{n_c\}) = \prod_c \binom{N_c}{n_c} \binom{N}{n}^{-1}.$$

In the case of the restricted AR sampling, the law is multinomial:

$$\Pr(\{n_c\}) = \prod_c p_c^{n_c} / n_c!.$$

The sampling plan model retained by the quota method described in paragraph 3 corresponds to constraints on these two schemes; which is equivalent to constraints on the SR and AR plans.

If we assume that  $N$  tends toward infinity, that  $f$  tends towards 0, and that  $n^* = fN$  tends toward infinity, then in the two plans, the law of the  $u_c = n^{*-1/2} (n_c - fN_c) = n^{*1/2} (p_c^* - p_c)$ , with  $p_c^* = n_c/n^*$ , tends toward a multidimensional normal law with independent  $u_c$ , with zero probability and variances equal to  $p_c$ .

### 3. Proportional Sampling

In this case, we have  $\hat{N}_c = N/n n_c$ , so that the quantity for which we want to determine the variance is:

$$\frac{N}{n^{*1/2}} \sum_c u_c E_c,$$

where the vector of the  $u_c$  follows a centered normal law with a diagonal covariance matrix  $\Delta = \text{diag}(p_c)$ , constrained by the relationships expressed by the quotas:

$$\sum_{q_c=i} u_c = 0 \text{ for } q = 1 \text{ to } Q, \quad i = 1 \text{ for } I_q \text{ if } q = 1, \quad i = 1 \text{ for } I_q - 1 \text{ if } q = 2 \text{ for } Q.$$

If we let  $U$  represent the vector of the  $u_c$ , the relationships can be written as follows:

$$AU = 0,$$

with  $A$  matrix with  $l = \sum_q I_q - (Q - 1)$  rows and  $k = \Pi_q I_q$  columns, where 1 and 0 represent the constraints. This also expresses the fact that  $U$  varies in the kernel  $L$  of the operator defined by matrix  $A$ . The (asymptotic) law of  $U$  is thus that of a centered gaussian vector  $W$  with a matrix whose covariances equal  $\Delta$ , when  $AW = 0$ . Thus, it is a question of evaluating the variance of a scalar product  $U'E$ , where  $\underline{E}$  is the vector of the  $E_c$ .

It is important to emphasize the following two points:

- The constraints upon the  $E_c$  given in result 1 can be expressed on the basis of matrix analysis by  $A\underline{\Delta E} = 0$ . In other words,  $\underline{\Delta E}$  is a vector of  $L = \text{Ker}A$ , or a vector of  $\text{Ker}(A\underline{\Delta})$ .
- Let  $P$  be the projection of  $\mathfrak{R}^k$  on  $L$  orthogonal in the  $\Delta^{-1}$  metrics.  $P$  verifies the following relations:
  - $\forall x \in L, Px = x; \text{Im } P = L$
  - $P y = 0 \Leftrightarrow \forall x \in L, x' \Delta^{-1} y = 0; \text{Ker } P = \Delta(L^\perp)$ ,

where  $L^\perp$  is the supplementary line orthogonal to  $L$  in the natural metrics.

The gaussian vectors  $PW$  and  $(1 - P)W$  vary in  $L$  and  $\Delta(L^\perp)$  respectively; and their sum is equal to  $W$ . Moreover, they are independent; in fact, their covariance matrix is  $E(PW)((1 - P)W)' = P\underline{\Delta}(1 - P')$ . Thus,  $P'$  is the kernel projector  $L^\perp$  and can be represented as  $\Delta(L^\perp)^\perp$ . The image of the projector  $(1 - P')$  is thus  $L^\perp$ . That of  $\Delta(1 - P')$  is  $\Delta(L^\perp)$ ; that is, the kernel of  $P$ , *q.e.d.*

At this point, we have to evaluate the variance of  $\sum_c u_c E_c = U'E$ . Thus, in accordance with the previous statements, we can write  $W = U + V$ , when  $U$  and  $V$  are independent. The law of  $W$  conditional upon  $W \in L$  is none other than the law of  $W$  conditional upon  $V = 0$ .

Moreover, we have:

$$V'E = (\Delta^{-1} V)' (\Delta E).$$

Since  $\Delta E$  is in  $L$ , and  $V$  varies in  $\Delta(L^\perp)$ , the scalar product above is zero. From this, we can deduce that:

$$\text{Var}(U'E) = \text{Var}(W'E) = \underline{E}' \underline{\Delta E} = \sum_c p_c E_c^2.$$

The asymptotic variance of is thus equal to  $N/n^* \sum_c n_c E_c$

$$\frac{N^2}{n} \sum_c p_c E_c^2 = \frac{N}{n} \sum_c N_c E_c^2.$$

#### 4. Sampling using "Non-Proportional" Quotas

Let us complete the preceding asymptotic reduction. Now, the vector  $\hat{p}^\circ$  of  $n_c/n^*$  is constrained by

$$A\hat{p}^\circ = Ap + n^{*-1/2} AV_0,$$

where  $Ap$  is the vector (1-dimensional) of the “proportional quotas”, and  $V_0$  is the only vector ( $k$ -dimensional) of  $\Delta(L^\perp)$ , so that  $A(p + n^{*-1/2} V_0)$ ; that is, the vector of the quotas imposed. Thus, as in the previous paragraph,  $U = n^{*1/2} (\hat{p}^\circ - p)$  may be analyzed as a gaussian vector  $W = U + V$  conditional upon  $V = V_0$ . Thus,  $EU_0 = V_0$ , and the covariances matrix of  $U_0$  is the same as that of  $U$ .

Moreover, we go from  $\hat{p}^\circ$  to  $\hat{p}$  by estimating the maximum resemblance. Under asymptotic gaussian conditions, this consists of minimizing the quadratic form  $(\hat{p}^\circ - \hat{p})' \Delta^{-1} (\hat{p}^\circ - \hat{p})$  under constraints  $A\hat{p} = Ap$ . Since  $\hat{p}^\circ$  varies in the related subspace  $L + V_0$  that is parallel to  $L$ , and minimization is a question of projecting  $\hat{p}^\circ$  upon  $L$  orthogonally for  $\Delta^{-1}$ ; that is, along  $\Delta(L^\perp)$ , it follows that we have  $\hat{p} = \hat{p}^\circ - n^{*-1/2} V_0$  under asymptotic conditions. The random vector  $\hat{p}$  is thus obtained from  $\hat{p}^\circ$ , is unbiased, and has the same covariance matrix as  $\hat{p}^\circ$ , so that  $n^{*-1/2} U$ .

Finally, we have:

$$E \left( \sum_c \hat{p}_c E_c \right)^2 = E(\hat{p}' \underline{E})^2 = \frac{1}{n^*} \sum_c p_c E_c^2$$

as in the previous case.

## REFERENCES

- CASSEL, C.M., SÄRNDAL, C.E., and WRETMAN, J.H. (1977). *Foundations of Inference in Survey Sampling*. New York: Wiley & Sons.
- COCHRAN, W.G. (1977). *Sampling Techniques*. New York: Wiley & Sons.
- DESABIE, J. (1965). *Théorie et pratique des sondages*. Paris: Dunod.
- DEVILLE, J.C., and SÄRNDAL, C.E. (1990). Calibration estimators and generalized raking techniques. Manuscript submitted for publication.
- GOURIÉROUX, C. (1981). *Théorie des sondages*. Paris: Economica.
- MADOW, W.G., OLKIN, I., and RUBIN, D.B., (Eds.) (1983). *Incomplete Data in Sample Surveys*. New York: Academic Press.
- RAO, J.N.K. (1976). Unbiased variance estimation for multistage designs. *Sankhyā*, Series C, 37, 133-139.
- SMITH, T.M.F. (1983). On the validity of inferences from non-random samples. *Journal of the Royal Statistical Society*, A, 146, 394-403.
- WOLTER, K.M. (1985). *Introduction to Variance Estimation*. New York: Springer-Verlag.