Ivan P. Fellegi Dominion Bureau of Statistics, Ottawa

(2)

# Introduction

Let the population consist of N units and assume that a "measure of size",  $m_i$  (i=1, 2, ..., N), is associated with each unit. A general problem of selecting n units with probability proportional to size (pps) without replacement is to devise sampling schemes such that  $\pi_i$ , the probability that the i-th unit is

in the sample is proportional to m<sub>i</sub>, i.e.

 $p_i = \frac{m_i}{\Sigma m_i}$ 

$$r_{i} = nm_{i}/\Sigma m_{i}$$
 (1)

Setting

(1) can be rewritten as 
$$\pi_i = np_i$$
 (3)

The present author proposed a scheme [2] of successively selecting n units without replacement in such a way that the probability  $\delta_i^k$  of selecting the i-th unit at the k-th selection equals  $p_i$ :

$$\delta_{i}^{k} = P_{i}$$
 i=1,2,...,N; k=1,2,...,n (4)

Since the events of selecting unit i at any two draws are mutually exclusive, it follows that

$$\pi_{i} = \sum_{k=1}^{n} \delta_{i}^{k} = np_{i}$$

This general method of selection has some particularly useful features for n=2. A method has been outlined in [2] (and a FORTRAN programme is available) to solve the following system of equations for  $q_i$  (i=1,2,...,N) and a:

$$aq_{i}(1-q_{i}) = 2p_{i}$$
  $i=1,2,...,N$  (5)

$$\sum_{i=1}^{N} q_i = 1$$
 (6)

$$a = 2/(1 - \sum_{i=1}^{N} q_i^2)$$
 (7)

The selection of two units is then carried out<sup>\*</sup> by selecting the first unit with probabilities  $\{p_i; i=1,2,\ldots,N\}$  and having selected unit i at

\* Brewer [1] proved that equations (5) to (7) always have a unique solution if only  $p_i < \frac{1}{2}$  (i=1,2,...,N). the first draw, selecting the second unit from among the remaining units with probabilities  $\{q_j/(1-q_i); j=1,2,\ldots,N; j\neq i\}$ . The probability of selecting units i and j in that order is, therefore, equal to

$$p_{i} \frac{q_{i}}{1-q_{i}}$$
(8)

while the probability of selecting the same units in reverse order is equal to

$$p_{j} \frac{q_{i}}{1-q_{i}}$$
 (9)

It follows from (5) that the expressions under (8) and (9) are both equal to

$$\frac{1}{2} aq_{i}q_{j} = p_{i} \frac{q_{j}}{1-q_{i}} = p_{j} \frac{q_{i}}{1-q_{j}}$$
(10)

and hence  $\pi_{ij}$ , the probability of selecting units i and j in either order is equal to

$$\pi_{ij} = aq_i q_j \tag{11}$$

Consequently, given that units i and j are in the sample, the conditional probability of selecting them in either order equals 1/2. Formula (4) follows from formulae (5) to (11).

The problem posed in the present paper is the following. Suppose a sample of two units had been selected as described above, with probabilities proportional to the measures of size {m<sub>i</sub>}. The two units might, for example, be two primary sampling units (areas) and the measures of size might be census counts. Some years later a new census is taken, resulting in new measures of size  $M_i$ . It is desirable to select a sample of two psu's with the following properties: I) The probabilities of selection are proportional to the new measures of size. II) Since, however, the original sample of two psu's often represents a capital investment (listings of households might have been prepared, enumerators trained, etc.), it is desirable to maximize the overlap between the old and new sample. More precisely, if X is a random variable (depending on the sampling scheme for the selection of the new sample), such that

- X = 0 if the new sample coincides with the old
  - = 1 if the new and old samples have one unit in common

# = 2 if the new and old samples are disjoint

then it is desirable to have a sampling scheme for the selection of the new sample for which

E(X) = expected number of rejections

is minimized. III) it is also desirable to have the joint probabilities of selecting two units in the new sample satisfy equations (4) to (11) in terms of  $P_i$ ,  $Q_i$ , A and  $\Pi_{ij}$  (defined in terms of the new size measures in analogy to  $P_i$ ,  $q_i$ , a and  $\pi_{ij}$ ) so that the process of revising the size measures in the future can be repeated again (actually property III implies I).

The procedure outlined below achieves I) and III) and achieves II) approximately.

# The procedure

Keyfitz [3] provided a procedure to achieve the objectives outlined above when the sample consists of one unit. In this simple case sampling with and without replacement become indistinguishable and property I) automatically implies III). Since, however, the present procedure follows the methods of Keyfitz, it is instructive to recapitulate them here briefly.

A sample of one unit had been selected with probabilities  $\{p_j\}$ . The desired new probabilities are  $\{P_j\}$ . The procedure of changing the probabilities of selection is summarized as follows: first the original unit is subjected to a test of retention; it is retained with probability

R<sub>j</sub>

rejected with probability

1-R

In case of rejection, the second step consists of a subsequent selection. The new probability of selecting unit j is therefore

$$P_{j} = P_{j}R_{j} + S_{j}$$

S being the probability of selecting unit j at the second step. The objective is to maximize  $R_j$  (or minimize S ) within the equation above. Clearly if

 $R_i < 1$  and  $S_i > 0$ 

the procedure is not optimal, since  $R_j$  could be further increased,  $S_j$  further decreased. In this case the probability of rejection is positive

$$1-R_{1} > 0$$

Hence a procedure is not optimal if it permits the rejection of a unit followed by its subsequent re-selection. The converse is also true: a procedure is optimal if the set of units with a positive probability of rejection and the set of units with a positive probability of subsequent selection are disjoint. The Keyfitz procedure is optimal in this sense. The set of units with positive probability of rejection is the set of decreasing units for which

 $P_j > P_j$ 

None of these units can be selected subsequently. The set of units with positive probability of subsequent selection is the set of units which increased, i.e. for which

<sup>p</sup>j ≦ <sup>P</sup>j

Formally

$$I = \{i: p_i \leq P_i\} \text{ and } D = \{d: p_d > P_d\}$$
(12)

The Keyfitz procedure consists of retaining the originally selected unit if it is in I. If the originally selected unit (d) is in D, then it is retained with probability  $R_d = P_d/p_d$ ,

a unit in D is rejected, then a unit in I is selected with probabilities

$$S_{i} = \frac{P_{i} - P_{i}}{\sum (P_{i} - P_{i})} \qquad i \in I \qquad (13)$$

It is easy to show that the probability of unit j being in the sample after the adjustment is equal to  $P_j$  (j=1,2,...,N) and that the procedure

is optimal in the sense of property II) above (it minimizes the expected number of rejections). These properties are based on the observation that since

$$\begin{matrix} N & N \\ \Sigma & P_j &= \Sigma & P_j &= 1 \\ j = 1 & j & j = 1 \end{matrix}$$

therefore

$$\sum_{d \in D} (p_d - P_d) = \sum_{i \in I} (P_i - P_i)$$
(14)

The present procedure starts with the observation that the sampling procedure by which the original two units were selected is symmetric [2] i.e. the conditional probability of the two units having been selected in either order is equal to 1/2. Consequently if one does not know the order in which the units had been selected, one can "recreate" the order at random, e.g. by selecting a random number r  $(0 < r \leq 1)$  and assuming that if  $0 < r \leq 0.5$  than the order in which units i and j had been selected was higher subscript first, lower subscript second; if  $0.5 < r \leq 1$  then the order in which the units had been selected.

After this step we may assume that we have an ordered sample of two units (i first, j second) whose joint probability of selection was

$$a_{j}^{\pi} = p_{i} \frac{q_{j}}{1-q_{i}} = \frac{1}{2} a q_{i} q_{j}$$
 (15)

Next we observe that the (unconditional) probability of unit i having been selected first was equal to p<sub>i</sub>. Since (4) is

to be maintained (in terms of the new measures) after the adjustment, the probability of unit i to be selected as the first should be equal to  $P_i$ . The Keyfitz procedure can be applied to

this end without modification: the two sets I and D, defined under (12), are formed and the Keyfitz rules for retention, rejection and reselection are applied.

After the adjustment of the probability of selection of the first unit the probability of selection of the second unit is adjusted. Here the Keyfitz principle has to be modified slightly. To illustrate this point, consider the sets

$$I_{i} = \{j: j \neq i \text{ and } \frac{q_{i}}{1-q_{i}} \leq \frac{Q_{i}}{1-Q_{i}}\} \qquad i=1,2,\ldots,N$$
$$D_{i} = \{j: j \neq i \text{ and } \frac{q_{i}}{1-q_{i}} > \frac{Q_{i}}{1-Q_{i}}\}$$

where  $q_j/(1-q_i)$  and  $Q_j/(1-Q_i)$  are the conditional probabilities of selecting unit j after i was selected.

# Suppose in (15) is D and $j \in D_{i}$ .

According to the Keyfitz principle unit i is retained with probability

 $\frac{P_i}{P_i}$ 

and unit j is retained with probability

$$\frac{Q_{i}}{1-Q_{i}} / \frac{q_{i}}{1-q_{i}}$$

This adjustment yields a probability of selection for units i and j in that order

$${}_{i}^{\Pi}{}_{j} = {}_{i}^{\pi}{}_{j} \frac{{}^{P}{}_{i}}{{}_{P}{}_{i}} \cdot \frac{{}^{Q}{}_{i}}{1 - {}_{Q}{}_{i}} / \frac{{}^{q}{}_{j}}{1 - {}_{q}{}_{i}}$$
$$= {}_{P}{}_{i} \frac{{}^{Q}{}_{i}}{1 - {}_{Q}{}_{i}} = \frac{1}{2} A {}_{Q}{}_{i} {}_{Q}{}_{j}$$
(16)

provided that neither unit has a chance of being re-selected once rejected (which is part of the Keyfitz principle). Formula (16) is the desirable equivalent of (15) after the adjustment. This reasonably simple procedure, however, has to be modified (as far as the adjustment of the probability of selection of the second unit is concerned!) if the first unit happens to be rejected. The reason for this is the fact that as long as the first unit is retained the Keyfitz procedure can be applied at the second step to the resulting conditional distribution (the condition being the selection of unit i at the first step). If, however, the first unit is rejected and another one is selected (e.g. i') then one has to compare

$$\frac{q_i}{1-q_i}$$

with

i.e. one is compelled to consider two different conditional distributions.

The following modification of the Keyfitz principle yields the desirable result:

#### Theorem 1.

Suppose that a sample of two units had been selected without replacement with probabilities proportional to the measures of size  $\{m_i\}$  satisfying formulae (1) to (11) (using the procedure described in [2]). Let  $\{M_i\}$  be the new measures of size and let  $\{P_i\}$ ,  $\{Q_i\}$  and A be computed from formulae (2), (5), (6) and (7). Then the application of the rules below will result in a sample having the properties (1) to (11) in terms  $\{M_i\}$ ,  $\{P_i\}$ ,  $\{Q_i\}$ , and A.

Define the subsets of the set of integers between 1 and N as follows

$$I = \{c: p_c \leq P_c\}$$
(17)

$$D = \{c: p_c > P_c\}$$
(18)

$$I_{j} = \{c: c \neq j \text{ and } \frac{q_{c}}{1-q_{j}} \leq \frac{Q_{c}}{1-Q_{j}}\}$$
for j=1,2,...,N (19)

$$D_{j} = \{c: c \neq j \text{ and } \frac{q_{c}}{1-q_{j}} > \frac{Q_{c}}{1-Q_{j}}\}$$
for j=1,2,...,N (20)

$$I_{id} = \{c: c \neq i, c \neq d \text{ and } \frac{q_c}{1 - q_i - q_d} \leq \frac{Q_c}{1 - Q_i} \}$$

$$D_{id} = \{c: c \neq i, c \neq d \text{ and } \frac{q_c}{1 - q_i - q_d} > \frac{Q_c}{1 - Q_i}\}$$
for icl and dcD (22)

Denote

$$S = \sum_{c \in I} (P_c - P_c) = \sum_{c \in D} (p_c - P_c)$$
(23)

$$S_{j} = \sum_{c \in I_{j}} \left( \frac{Q_{c}}{1 - Q_{j}} - \frac{q_{c}}{1 - q_{j}} \right)$$
$$= \sum_{c \in D_{j}} \left( \frac{q_{c}}{1 - q_{j}} - \frac{Q_{c}}{1 - Q_{j}} \right)$$
(24)

$$S_{id} = \sum_{c \in I_{id}} \left( \frac{Q_c}{1 - Q_i} - \frac{q_c}{1 - q_i^{-q_d}} \right)$$
(25)

Rule 1: Given the original sample of two units determine the order in which they will be subjected to a test of rejection by selecting a random number r ( $0 < r \le 1$ ). If  $0 < r \le 0.5$ the unit with the lower subscript will be adjusted first (called "the first unit"); if  $0.5 < r \le 1$  then the unit with the higher subscript will be the "first unit".

Rule 2: Apply to the first unit the Keyfitz rule, i.e. a) if unit c is the first and ceI then retain it; b) if ceD retain it with probability  $P_c/p_c$ , reject with probability  $1-P_c/p_c$  and in this case select a unit from I with probability  $(P_c-P_c)/S$ .

Rule 3: If the first unit was retained (with subscript j) and the second unit had subscript k:

- (a) retain the second unit if kel
- (b) if kcD, retain it with probability

$$\frac{\mathsf{Q}_{\mathsf{k}}}{\mathsf{1}-\mathsf{Q}_{\mathsf{j}}} \; / \; \frac{\mathsf{q}_{\mathsf{k}}}{\mathsf{1}-\mathsf{q}_{\mathsf{j}}} = \frac{\mathsf{Q}_{\mathsf{k}}(\mathsf{1}-\mathsf{q}_{\mathsf{j}})}{\mathsf{q}_{\mathsf{k}}(\mathsf{1}-\mathsf{Q}_{\mathsf{j}})}$$

reject it otherwise; if it is rejected, select another one from I<sub>i</sub> with probabilities

$$\frac{1}{S_j} \left( \frac{Q_c}{1-Q_j} - \frac{q_c}{1-q_j} \right); \qquad c \epsilon I_j$$

Rule 4: If the first unit (d) was in D and if it was rejected and replaced by  $i \in I$  and if the second unit had subscript k

- (a) retain the second unit if kelid
- (b) if kcD<sub>id</sub>, retain it with probability

$$\frac{Q_k}{1-Q_i} / \frac{q_k}{1-q_i-q_d} = \frac{Q_k(1-q_i-q_d)}{q_k(1-Q_i)}$$

reject it otherwise; if it is rejected, select d with probability 0.

$$\frac{d}{Q_d + (1 - Q_i) S_{id}}$$

and if d is not selected, select a unit from  $I_{id}$  with probabilities

$$\frac{1}{S_{id}} \left( \frac{Q_c}{1-Q_i} - \frac{q_c}{1-q_i-q_d} \right), \qquad c \epsilon I_{id}$$

(c) if k=i, select from among all the units, excepting i, with probabilities

$$\frac{Q_{c}}{1-Q_{i}} \qquad c=1,2,\ldots,N; \quad c\neq i$$

Two lemmas are required for the proof of theorem 1:

Lemma 1. Let S<sub>id</sub> be defined by formula (25). Then

$$S_{id} = \sum_{c \in I_{id}} \left( \frac{Q_c}{1 - Q_i} - \frac{q_c}{1 - q_i - q_d} \right)$$
$$= \sum_{c \in D_{id}} \left( \frac{q_c}{1 - q_i - q_d} - \frac{Q_c}{1 - Q_i} \right) - \frac{Q_d}{1 - Q_i} \quad (26)$$

Lemma 2. Suppose that in the original sample the first unit (as obtained from Rule 1) was unit dED, the second unit was  $k(\neq d)$  and that the application of Rule 2 resulted in the replacement of d by iEI. Denote the conditional probability of obtaining (through retention or through rejection and a new selection) unit  $\beta$  as the second unit of the new sample by

$$\sum_{\substack{k=1\\k\neq d}}^{N} q_k^P (d \neq i, k) = (1-q_d) \frac{Q_\beta}{1-Q_i}$$
(27)

for any dcD, icI and any  $\beta(\neq i)$ .

The proof of lemma 1 follows immediately from (6) (i.e. the analogous result for  $Q_i$ ) and from the observation that the union of  $D_{id}$  and  $I_{id}$  is the set of integers from 1 to N excepting i and d. The proof of lemma 2 will be presented after the proof of theorem 1.

#### Proof of theorem 1:

Denoting the first and second units of the new sample by  $\alpha$  and  $\beta$  respectively, there are four possibilities:

- (1) αεΙ, βεΙα
- (ii) αεΙ, βεD
- (111) αεD, βεΙ
- (iv) αεD, βεD

438

It will be shown that in each case the probability of obtaining units  $\alpha$  and  $\beta$  in that order in the new sample is equal to

$$\frac{1}{2} AQ_{\alpha}Q_{\beta}$$
 (28)

This will complete the proof of the theorem in that it will then follow from (5) (in terms of A,  $P_i$ ,  $Q_i$ ) that the probability of obtaining

in the new sample unit  $\alpha$  as the first unit is equal to

$$\sum_{\beta(\neq\alpha)} \frac{1}{2} AQ_{\alpha}Q_{\beta} = \frac{1}{2} AQ_{\alpha}(1-Q_{\alpha}) = P_{\alpha}$$
  
  $\alpha=1, 2, ..., N$ 

Also the probability of obtaining unit  $\boldsymbol{\beta}$  as the second unit is equal to

-

$$\sum_{\alpha \neq \beta} \frac{1}{2} AQ_{\alpha}Q_{\beta} = \frac{1}{2} A(1-Q_{\beta}) Q_{\beta} = P_{\beta}$$
  
 
$$\beta=1, 2, ..., N$$

and from the symmetry of (28) the probability of obtaining  $\alpha$  and  $\beta$  in either order is equal to

$$\Pi_{\alpha\beta} = AQ_{\alpha}Q_{\beta}$$

In the following proofs we assume that rule 1 had already been applied and that consequently the original sample had already been ordered.

Proof of case (i): 
$$\alpha \in I$$
,  $\beta \in I$ 

Such a sample can arise in the following three ways: the original sample was  $\alpha$  and  $\beta$  in that order in which case they are both retained (rules 2 and 3a); the original sample was  $\alpha$  and  $k\epsilon D_{\alpha}$ ,  $\alpha$  was retained and  $\beta\epsilon I_{\alpha}$  selected (rules 2 and 3b); the first unit was dcD, it was rejected,  $\alpha\epsilon I$  selected (rule 2) then  $\beta$  was obtained for the second unit as in lemma 2. The three probabilities corresponding to these three events are as follows:

$$\frac{1}{2} \operatorname{aq}_{\alpha} q_{\beta} + \frac{1}{2} \operatorname{aq}_{\alpha} q_{\beta} + \frac{1}{2} \operatorname{aq}_{\alpha} q_{k} \left[1 - \frac{Q_{k}(1 - q_{\alpha})}{q_{k}(1 - Q_{\alpha})}\right] \frac{1}{S_{\alpha}} \left(\frac{Q_{\beta}}{1 - Q_{\alpha}} - \frac{q_{\beta}}{1 - q_{\alpha}}\right) + \frac{\Sigma}{d \epsilon D} \sum_{k \neq d} \frac{1}{2} \operatorname{aq}_{d} q_{k} \left(1 - \frac{P_{d}}{P_{d}}\right) \frac{1}{S} (P_{\alpha} - P_{\alpha}) P_{\beta} (d \neq \alpha, k) = \frac{1}{2} \operatorname{aq}_{\alpha} q_{\beta} + \frac{1}{2} \operatorname{aq}_{\alpha} (1 - q_{\alpha}) \frac{1}{S_{\alpha}} \left(\frac{Q_{\beta}}{1 - Q_{\alpha}} - \frac{q_{\beta}}{1 - q_{\alpha}}\right) \sum_{k \in D_{\alpha}} \left(\frac{q_{k}}{1 - q_{\alpha}} - \frac{Q_{k}}{1 - Q_{\alpha}}\right) + \frac{1}{S} \left(P_{\alpha} - P_{\alpha}\right) \sum_{d \epsilon D} \frac{1}{2} \operatorname{aq}_{d} (1 - \frac{P_{d}}{P_{d}}) \left(1 - q_{d}\right) \frac{Q_{\beta}}{1 - Q_{\alpha}}$$

the last term following from lemma 2. Applying (5), (10), (23) and (24) we obtain

$$\frac{1}{2} aq_{\alpha}q_{\beta} + \frac{1}{1-Q_{\alpha}} - \frac{1}{2} aq_{\alpha}q_{\beta} + p_{\alpha} \frac{Q_{\beta}}{1-Q_{\alpha}} - \frac{1}{2} aq_{\alpha}q_{\beta} + (P_{\alpha}-P_{\alpha}) \frac{Q_{\beta}}{1-Q_{\alpha}} = P_{\alpha} \frac{Q_{\beta}}{1-Q_{\alpha}} = \frac{1}{2} AQ_{\alpha}Q_{\beta}$$
Proof of case (11):  $\alpha \in I$ ,  $\beta \in D_{\alpha}$ 

Such a sample can arise in two ways: the original sample was  $\alpha$  and  $\beta$  and  $\beta$  was retained ( $\alpha$  is certainly retained by rule 2); the first unit was dcD, it was rejected,  $\alpha$ cI selected (rule 2) then  $\beta$  was obtained for the second unit as in lemma 2. The probabilities corresponding to these events are as follow:

$$\frac{1}{2} aq_{\alpha}q_{\beta} \frac{Q_{\beta}(1-q_{\alpha})}{q_{\beta}(1-Q_{\alpha})} +$$

$$+ \sum_{d \in D} \sum_{k(\neq d)} \frac{1}{2} aq_{d}q_{k} (1-\frac{P_{d}}{P_{d}}) \frac{1}{S}(P_{\alpha}-P_{\alpha}) P_{\beta} (d+\alpha, k)$$

$$= \frac{1}{2} aq_{\alpha}(1-q_{\alpha}) \frac{Q_{\beta}}{1-Q_{\alpha}} +$$

$$+ (P_{\alpha}-P_{\alpha}) \frac{Q_{\beta}}{1-Q_{\alpha}}$$

the second term on the right hand side following from lemma 2 in a way identical to the manipulation of the last term in the proof of (i). Applying (5) to the first term we obtain

$$p_{\alpha} \frac{Q_{\beta}}{1-Q_{\alpha}} + (P_{\alpha}-P_{\alpha}) \frac{Q_{\beta}}{1-Q_{\alpha}} = P_{\alpha} \frac{Q_{\beta}}{1-Q_{\alpha}} = \frac{1}{2} AQ_{\alpha}Q_{\beta}$$
Proof of case (iii):  $\alpha \in D$ ,  $\beta \in I_{\alpha}$ 

Such a sample can arise in two ways: the original sample was  $\alpha$  and  $\beta$  and  $\alpha$  was retained (rule 2;  $\beta$  is then automatically retained by rule 3a); the original sample was  $\alpha$  and  $k\epsilon D_{\alpha}$ ,  $\alpha$  retained (rule 2) k rejected and  $\beta\epsilon I_{\alpha}$  selected (rule 3b). The probabilities corresponding to these events are as follow:

$$\frac{1}{2} aq_{\alpha}q_{\beta} \frac{P_{\alpha}}{P_{\alpha}} + \frac{\Sigma}{k \varepsilon D_{\alpha}} \frac{1}{2} aq_{\alpha}q_{k} \frac{P_{\alpha}}{P_{\alpha}} (1 - \frac{Q_{k}(1-q_{\alpha})}{q_{k}(1-Q_{\alpha})}) \frac{1}{S_{j}} (\frac{Q_{\beta}}{1-Q_{\alpha}} - \frac{q_{\beta}}{1-q_{\alpha}}) =$$

$$= \frac{1}{2} aq_{\alpha}q_{\beta} \frac{P_{\alpha}}{P_{\alpha}} +$$

$$+ \frac{1}{2} aq_{\alpha}(1-q_{\alpha}) \frac{P_{\alpha}}{P_{\alpha}} \frac{1}{S_{j}} (\frac{Q_{\beta}}{1-Q_{\alpha}} - \frac{q_{\beta}}{1-q_{\alpha}})$$

$$\sum_{k \in D_{\alpha}} (\frac{q_{k}}{1-q_{\alpha}} - \frac{Q_{k}}{1-Q_{\alpha}})$$

$$= P_{\alpha} \frac{q_{\beta}}{1-q_{\alpha}} +$$

$$+ P_{\alpha} (\frac{Q_{\beta}}{1-Q_{\alpha}} - \frac{q_{\beta}}{1-q_{\alpha}})$$

the first term following from (5), the second from (5) and (24). We obtain

$$P_{\alpha} \frac{Q_{\beta}}{1-Q_{\alpha}} = \frac{1}{2} A Q_{\alpha} Q_{\beta}$$

Proof of case (iv):  $\alpha \in D$ ,  $\beta \in D_{\alpha}$ 

Such a sample can only arise in one way: the original sample was  $\alpha$  and  $\beta$ ,  $\alpha$  was retained (rule 2) and  $\beta$  was retained (rule 3b). The probability of this event is

$$\frac{1}{2} \operatorname{aq}_{\alpha} q_{\beta} \frac{P_{\alpha}}{P_{\alpha}} \frac{Q_{\beta}(1-q_{\alpha})}{q_{\beta}(1-Q_{\alpha})} = P_{\alpha} \frac{Q_{\beta}}{1-Q_{\alpha}} = \frac{1}{2} \operatorname{AQ}_{\alpha} Q_{\beta}$$

after two applications of (5).

1.

This completes the proof of theorem

Proof of lemma 2:

The original sample consisted of units dcD and k( $\neq$ d). The application of rule 2 resulted in the rejection of d and its replacement by icI. Lemma 2 is to be proved for all  $\beta(\neq i)$ , i.e. for the cases when  $\beta \epsilon D_{id}$ ,  $\beta \epsilon I_{id}$  and  $\beta$ =d.

(a) If  $\beta \in D_{id}$ :

According to rule 4a) and rule 4b) if  $k\epsilon I_{id}$  or if  $k\epsilon D_{id}$  but  $k\neq\beta$ 

$$P_o(d \rightarrow i, k) = 0$$

According to rule 4b) if  $k=\beta \in D_{id}$ 

$$P_{\beta} (d \neq i, \beta) = \frac{Q_{\beta}(1-q_i-q_d)}{q_{\beta}(1-Q_i)}$$

According to rule 4c)

$$P_{\beta} (d \neq i, i) = \frac{Q_{\beta}}{1-Q_{i}}$$

Consequently

$$\sum_{\substack{k=1\\k\neq d}}^{N} q_k P_\beta(d + i, k) = q_\beta \frac{Q_\beta(1 - q_i - q_d)}{q_\beta(1 - Q_i)} + q_i \frac{Q_\beta}{1 - Q_i}$$
$$= (1 - q_d) \frac{Q_\beta}{1 - Q_i}$$

(b) If Belid

According to rule 4a) if  $k \in I_{id}$  but  $k \neq \beta$  then

$$P_{\beta}(d \rightarrow i, k) = 0$$

and

$$P_{\beta}(d \rightarrow i, \beta) = 1$$

According to rule 4c)

$$P_{\beta}(d \neq i, i) = \frac{Q_{\beta}}{1-Q_{i}}$$

According to 4b) if kED id

$$P_{\beta}(d \neq i, k) = (1 - \frac{Q_{k}(1 - q_{i} - q_{d})}{q_{k}(1 - Q_{i})})(1 - \frac{Q_{d}}{Q_{d} + (1 - Q_{i})S_{id}})$$
$$\frac{1}{S_{id}}(\frac{Q_{\beta}}{1 - Q_{i}} - \frac{q_{\beta}}{1 - q_{i} - q_{d}})$$

Hence

$$\sum_{\substack{k=1\\k=1}}^{N} q_{k} P_{\beta}(d \neq i, k) = q_{\beta} + q_{i} \frac{Q_{\beta}}{1 - Q_{i}} + (1 - \frac{Q_{d}}{Q_{d} + (1 - Q_{i})S_{id}}) \frac{1}{S_{id}} (\frac{Q_{\beta}}{1 - Q_{i}} - \frac{q_{\beta}}{1 - q_{i} - q_{d}})$$

$$(1 - q_{i} - q_{d}) \sum_{\substack{k \in D_{id}}} (\frac{q_{k}}{1 - q_{i} - q_{d}} - \frac{Q_{k}}{1 - Q_{i}})$$

Applying lemma 1 to the last term we obtain

$$q_{\beta} + q_{i} \frac{Q_{\beta}}{1-Q_{i}} + \frac{(1-Q_{i})S_{id}}{Q_{d}+(1-Q_{i})S_{id}} \frac{1}{S_{id}} (\frac{Q_{\beta}}{1-Q_{i}} - \frac{q_{\beta}}{1-q_{i}-q_{d}})$$

$$(1-q_{i}-q_{d}) (S_{id} + \frac{Q_{d}}{1-Q_{i}})$$

$$= (1-q_{d}) \frac{Q_{\beta}}{1-Q_{i}}$$

$$(c) \text{ If } \beta=d$$

According to rule 4b) if kEI<sub>18</sub>

$$P_{\beta} (\beta \rightarrow i, k) = 0$$

According to rule 4c)

$$P_{\beta} (\beta \neq i, i) = \frac{Q_{\beta}}{1 - Q_{i}}$$

According to rule 4b) if kcD<sub>i8</sub>

$$P_{\beta}(\beta \rightarrow i, k) = (1 - \frac{Q_{k}(1 - q_{i} - q_{\beta})}{q_{k}(1 - Q_{i})}) \frac{Q_{\beta}}{Q_{\beta} + (1 - Q_{i})S_{i\beta}}$$

Consequently

$$\sum_{\substack{k=1\\(k\neq\beta)}}^{N} q_k P_{\beta}(\beta+i,k) = q_i \frac{Q_{\beta}}{1-Q_i} + \frac{Q_{\beta}}{Q_{\beta}+(1-Q_i)S_{i\beta}} (1-q_i-q_{\beta})$$

$$\sum_{\substack{k\in D_{i\beta}}} (\frac{q_k}{1-q_i-q_{\beta}} - \frac{Q_k}{1-Q_i})$$

Applying lemma 1 to the last term, we obtain

$$q_{i} \frac{Q_{\beta}}{1-Q_{i}} + \frac{Q_{\beta}}{Q_{\beta}+(1-Q_{i})S_{i\beta}} (1-q_{i}-q_{\beta}) (S_{i\beta} + \frac{Q_{\beta}}{1-Q_{i}})$$
$$= (1-q_{\beta}) \frac{Q_{\beta}}{1-Q_{i}}$$

This completes the proof of Lemma 2.

## The expected number of rejections

It has been proved in the previous section that the procedure as described by the rule of theorem 1 satisfies objectives I) and III) as set out in the introduction. With respect to the objective of minimizing the expected rejections the following may be said. The procedure would be optimal if it rejected originally selected units only to the extent necessary, i.e. if it rejected a unit only if its probability of selection for the new sample had to be diminished, and if in this case the original probability of selection times the probability of retention were equal the desired new probability of rejection. A procedure is not optimum if a rejected unit can be reselected since this means that the rejection procedure renders its probability for the new sample too small (it is rejected with an unnecessarily large probability).

The present procedure is such that a first unit (or a second unit), once rejected, can never be re-selected as a first unit (or a second unit). There is, however, one (and only one) condition, under which a unit rejected as a <u>first</u> unit can be re-selected as a <u>second</u> unit. This is embodied in rules 4b) and 4c) when a first unit dcD is rejected, icI is selected, the second unit k is in  $D_{id}$  or k=i, it is rejected and d is re-selected as a second unit. It is not too difficult to show that if one did not permit the re-selection of d as a

second unit at all, then the procedure would not yield the required probabilities for the new sample. Whether the present procedure is actually optimal, however, is not known. Neither is it known whether there is an optimal procedure at all (subject to objectives I) and III)). Yet the probability of the event described above which may result in a departure from optimality is sufficiently small to render the statement concerning the "near optimality" of the procedure more or less justified.

The actual formula for the expected number of rejections can be derived without difficulty.

# Theorem 2.

The expected number of rejections when applying the rules of theorem 1 is given by

$$\sum_{i \in I} P_{i} S_{i}^{+} \sum_{d \in D} P_{d} S_{d}^{+} S_{i}^{+} \sum_{i \in I} \sum_{d \in D} (P_{i}^{-} P_{i}^{-}) (P_{d}^{-} P_{d}^{-}) \frac{S_{id}^{-}}{S}$$

$$- \sum_{d \in D} \sum_{i \in I} \frac{1}{2} aq_{d} q_{i}^{-} (1 - \frac{P_{d}}{P_{d}^{-}}) \frac{P_{i}^{-} P_{i}^{-}}{S} \frac{Q_{d}^{-}}{1 - Q_{i}^{-}}$$

$$- \sum_{d \in D} \sum_{i \in I} \frac{1}{2} aq_{d} q_{i}^{-} (1 - \frac{P_{d}^{-}}{P_{d}^{-}}) \frac{P_{i}^{-} P_{i}^{-}}{S} S_{id}^{-} (29)$$

and an upper bound of (29) is given by

$$\sum_{i=1}^{N} P_i S_i + S$$
(30)

Proof:

The event "at least one rejection" is composed of the following mutually exclusive events:

- (i) the first unit is rejected
- (ii) the first unit is retained and the second unit rejected.

It is easy to show that the probability of event (i), as a result of applying rule 2, is given by

$$S = \sum_{i \in I} (P_i - P_i) = \sum_{d \in D} (P_d - P_d)$$
(31)

The probability of event (ii) is given by (rule 3)

$$\sum_{i \in I} \sum_{k \in D_{i}} \frac{1}{2} a q_{i} q_{k} (1 - \frac{Q_{k}(1 - q_{i})}{q_{k}(1 - Q_{i})}) + \sum_{d \in D} \sum_{k \in D_{i}} \frac{1}{2} a q_{d} q_{k} \frac{P_{d}}{P_{d}}$$

$$(1 - \frac{Q_{k}(1 - q_{d})}{q_{k}(1 - Q_{d})})$$

$$= \sum_{i \in I}^{\Sigma} p_i \sum_{k \in D_i}^{\Sigma} \left( \frac{q_k}{1 - q_i} - \frac{Q_k}{1 - Q_i} \right) + \sum_{d \in D}^{\Sigma} p_d \sum_{k \in D_i}^{\Sigma} \left( \frac{q_k}{1 - q_d} - \frac{Q_k}{1 - Q_d} \right)$$

$$= \sum_{i \in I} p_i S_i + \sum_{d \in D} P_d S_d$$
(32)

the middle line following from (5).

(31) plus (32) represent the probability of at least one rejection in the course of selecting the new sample. It is possible to arrive at the new sample via rejections and obtain the old sample. This happens according to rule 4c) with the probability (dcD and icI constituted the first sample in that order, d was rejected and i selected for the first unit in the new sample and d was then selected for the second unit in the new sample):

$$\sum_{\substack{d \in D \ i \in I}} \sum_{\substack{1 \ 2}} \frac{1}{2} a q_d q_i \quad (1 - \frac{P_d}{P_d}) \quad \frac{P_i - P_i}{S} \frac{Q_d}{1 - Q_i}$$
(33)

Consequently (31) plus (32) minus (33) represents the probability that the new sample is not identical with the old. Adding to this the probability of the event that the new sample will be entirely different from the old we shall obtain the expected number of rejections. The probability of the latter event (rule 4b) is given by

$$\sum_{d \in D} \sum_{i \in I} \sum_{k \in D_{id}} \frac{1}{2} aq_{d}q_{k} \left(1 - \frac{P_{d}}{P_{d}}\right) \frac{P_{1} - P_{1}}{S}$$

$$\left(1 - \frac{Q_{k}(1 - q_{1} - q_{d})}{q_{k}(1 - Q_{1})}\right) \left(1 - \frac{Q_{d}}{Q_{d} + (1 - Q_{1})S_{1d}}\right)$$

$$= \sum_{d \in D} \sum_{i \in I} \frac{1}{2}aq_{d} \left(1 - \frac{P_{d}}{P_{d}}\right) \frac{P_{1} - P_{1}}{S} \left(1 - q_{1} - q_{d}\right)S_{1d}$$

$$= \sum_{d \in D} \sum_{i \in I} \left(p_{d} - P_{d}\right) \left(P_{1} - p_{1}\right) \frac{S_{1d}}{S} - \sum_{d \in D} \sum_{i \in I} \frac{1}{2}aq_{d}q_{1}$$

$$(1 - \frac{P_d}{P_d}) \frac{P_i - P_i}{S} S_{id}$$
 (34)

the second line following from lemma 1.

Since the expected number of rejections is given by

$$(31)+(32)-(33)+(34)$$

this completes the proof of (29).

In order to prove (30), we note that  $I_{id}$  is a subset of  $I_i$ , consequently

hence

$$\sum_{i \in I}^{\Sigma} \sum_{d \in D} (P_i - P_i) (P_d - P_d) \frac{S_{id}}{S} \leq \sum_{i \in I}^{S} \sum_{d \in D} (P_i - P_d) \frac{S_{id}}{S} \leq \sum_{i \in I}^{S} \sum_{d \in D} (P_i - P_i) (P_i - P_d) \frac{S_{id}}{S} \leq \sum_{i \in I}^{S} \sum_{d \in D} (P_i - P_i) (P_i - P_d) \frac{S_{id}}{S} \leq \sum_{i \in I}^{S} \sum_{d \in D} (P_i - P_i) (P_i - P_d) \frac{S_{id}}{S} \leq \sum_{i \in I}^{S} \sum_{d \in D} (P_i - P_i) (P_i - P_d) \frac{S_{id}}{S} \leq \sum_{i \in I}^{S} \sum_{d \in D} (P_i - P_i) (P_i - P_i$$

$$\leq \sum_{i \in I}^{\Sigma} (P_i - P_i) \frac{P_d - P_d}{S} S_i$$
$$= \sum_{i \in I}^{\Sigma} (P_i - P_i) S_i$$
(35)

Omitting the last two terms in (29) and substituting the right hand side of (35) for the fourth term in (29) proves (30).

### An alternative procedure

An alternative procedure, also satisfying objectives I) and III) of the introduction is a direct application of the Keyfitz procedure to the entire sample (instead of the individual units). The procedure can be briefly described as follows:

Let  $\pi_{ij}$  be the probabilities of the original samples,  $\Pi_{ij}$  the desired new probabilities given by (11) in terms of A,  $Q_i$ ,  $Q_j$ . Define the sets of distinct samples (ordered pairs of integers) as follows:

$$K = \{(i,j): i < j \text{ and } \pi_{ij} \leq \Pi_{ij} \}$$
  
L =  $\{(i,j): i < j \text{ and } \pi_{ij} > \Pi_{ij} \}$ 

Let

$$\Gamma = \sum_{\substack{(i,j)\in K}} (\pi_{ij} - \pi_{ij}) = \sum_{\substack{(i,j)\in L}} (\pi_{ij} - \pi_{ij})$$

If the units of the old sample are in K, retain them. If they are in L retain them with probability  $\Pi_{ij}/\pi_{ij}$ . If the old sample is rejected select a new sample from K with probabilities

$$\frac{\pi_{ij}^{-\pi}_{ij}}{T}$$

Keyfitz's proof concerning his procedure applies here without change to show that this procedure satisfies objective III) and hence also I). It also follows from the Keyfitz proof that this procedure maximizes the probability of maintaining the entire old sample (i.e. it minimizes the probability of at least one rejection). In case of a rejection, however, the entire old sample is rejected. The procedure is not optimal, since it permits the re-selection of a rejected unit (if a sample is rejected and a new sample is selected, one of the units of the new sample may be identical to one of the units of the old).

Numerical examples show the relative performance of the two procedures: the alternative procedure results in 25 to 40 per cent (in one example 100 per cent) more rejections than the procedure of theorem 1. It is also easier to apply (although the rules appear to be more complicated) since under that procedure in any given concrete case at most 2N comparisons need to be made between old and new probabilities while the alternative procedure requires N(N-1)/2 such comparisons.

## Numerical examples

In examples 1 to 4 the same "original measures of size" are used.

In example 1 the "new measures of size" are obtained by reversing the ordering of the original measures of size: the unit with the smallest original size measure becomes the largest, etc. This represents a drastic change in the original size distribution. In example 2 the new size measures are obtained by interchanging consecutive pairs of the original size measures. This represents "intermittent moderate growth". In example 3 the fourth and fifth original measures of size are interchanged to obtain the new measures of size. This represents "some small growth and some small decline among the larger units". In example 4 the new measures are identical to the old except that the largest measure of size was doubled. This represents "large growth of a few units".

Examples 5 to 8 are analogous to examples 1 to 4 except that the population is larger (N=20) and there is a substantially higher variation in the original size measures.

The expected number of rejections, as in (29), and its upper bound, as in (30), is shown. Also the expected number of rejections under the alternative procedure of the previous section is shown.

	•		M <sub>i</sub>				
	1	<sup>m</sup> i	Ex. 1	Ex. 2	Ex. 3	Ex. 4	
	1	10	22	14	10	10	
	2	14	19	10	14	14	
	3	17	18	18	1/	1/	
	4	18	17	1/	19	10	
	5	19	14	22	18	19	
	6	22	10	19		44	
Expected No. of rejections		0.3759	0.1671	0.0239	0.3143		
Upper bound of Expected No. of rejections			0.4000	0.1752	0.0245	0.3237	
Alternative Expected No. of rejections			0.4998	0.1999	0.0342	0.4455	

Examples 1 to 4

Examples 5 to 8

	<sup>m</sup> i	Mi				
L 		Ex. 5	Ex. 6	Ex. 7	Ex. 8	
1	5	95	10	5	5	
2	10	90	5	10	10	
3	15	85	20	15	15	
4	20	80	15	20	20	
5	25	75	30	25	25	
6	30	70	25	30	30	
7	35	65	40	35	35	
8	40	60	35	40	40	
9	45	55	50	45	45	
10	50	50	45	50	50	
11	50	50	22	50	50	
12	22	45	50	22	22	
13	45	40	60	60	60	
14	70	30	75	70	70	
15	75	25	70	80	170	
10	80	20	85	75	80	
18	85	15	80	85	135	
10	90	10	95	90	90	
20	95	5	90	95	200	
Expected No. of rejections		0.9039	0.1012	0.0108	0.3066	
Upper bound Expected N of rejecti	of o. ons	0.9194	0.1027	0.0109	0.3109	
Alternative Expected No. of rejections		1.2170	0.1252	0.0193	0.4051	

In seven of the eight examples the alternative procedure yields between 20 to 45 per cent more rejections than the procedure of theorem 1. In example 7 the alternative procedure yields almost 75 per cent more rejections. In every case the upper bound for the expected number of rejections provided satisfactory approximations.

# References

- [1] Brewer, K.R.W.: "A note on Fellegi's method of sampling without replacement with probability proportional to size". To be published in the Journal of the American Statistical Association.
- [2] Fellegi, I.P.: "Sampling with varying probabilities without replacement: rotating and non-rotating samples". Journal of the American Statistical Association, 58 (1963), 183-201.
- [3] Keyfitz, N.: "Sampling with probabilities proportional to size: adjustment for changes in probabilities". Journal of the American Statistical Association, 46 (1951), 105-109.