

34094

UNITED NATIONS ECONOMIC AND SOCIAL COUNCIL



Distr.
LIMITED



E/CN.14/SM/3
23 February 1968

Original: ENGLISH

ECONOMIC COMMISSION FOR AFRICA
Seminar on Sampling Methods
Addis Ababa, 3 - 14 June 1968

SAMPLING FOR DEMOGRAPHIC AND HOUSING SURVEYS AND CIVIL REGISTRATION

M68-499

TABLE OF CONTENTS

	<u>Page</u>
1. INTRODUCTION	1- 3
2. GENERAL PRINCIPLES	
2.1 Sampling units and sampling frames.....	4- 5
2.2 Multi-stage sampling and clustering	5-10
2.3 Fixed versus variable sampling probabilities.....	10-14
2.4 Stratification.....	14-17
2.5 Sample size.....	17
2.6 Sampling in time.....	18
2.7 Relation to other sampling operations.....	18-19
3. SAMPLING OF RURAL SEDENTARY POPULATIONS FOR DEMOGRAPHIC CHARACTERISTICS	
3.1 Sampling units and sampling frames	
3.1.1 Area sampling units.....	20-21
3.1.2 Sampling of housing and households.....	22
3.2 Multi-stage sampling and clustering	
3.2.1 Variance between and within PSUs.....	23-24
3.2.2 Cost parameters.....	24-26
3.2.3 Cluster sampling.....	26-29
3.3 Details of sample design.	
3.3.1 Large PSUs to be split into smaller area units...	30-33
3.3.2 Large PSUs to be treated by sampling houses or compounds	33-34
3.4 Stratification.....	34
3.5 Sample size.....	34
3.6 Sampling in time	
3.6.1 Concentration or dispersion of the period of field work	35-36
3.6.2 Retrospective versus follow-up survey.....	36
3.6.3 Long-term retrospective data.....	37
3.7 Relation to other sampling operations.....	37

TABLE OF CONTENTS (cont'd)

	<u>Page</u>
4. SAMPLING OF URBAN POPULATION FOR DEMOGRAPHIC AND HOUSING CHARACTERISTICS	
4.1 Sampling units and sampling frames	
4.1.1 Area units.....	38
4.1.2 Units based on housing and property.....	38-39
4.1.3 Household lists, tax lists, etc.....	39
4.2 Clustering and multi-stage sampling.....	40
4.2.1 Cities or districts of regular lay-out.....	41
4.2.2 Cities or districts of irregular lay-out.....	41-44
4.2.3 Note on housing units.....	44
4.2.4 Note on housing surveys.....	45
4.2.5 Summary.....	45
4.3 Fixed versus variable sampling fractions.....	46
4.4 Stratification.....	46-47
4.5 Sample size.....	47
4.6 Sampling in time.....	47-49
4.7 Relation to other sampling operations.....	49-50
5. CIVIL REGISTRATION ON A SAMPLE BASIS FOR RURAL POPULATIONS	
5.1 General.....	51-52
5.2 Type of sample required.....	52-53
5.3 Base population.....	53-54

SAMPLING FOR DEMOGRAPHIC AND HOUSING SURVEYS AND CIVIL REGISTRATION

1. INTRODUCTION

This paper deals with three distinct problems:

1. Sampling of rural sedentary populations for demographic characteristics.
2. Sampling of urban populations for demographic and housing characteristics.
3. Civil registration on a sampling basis for rural populations.

Sampling methods are examined in the paper only in relation to the collection of data, although sampling can also be used at the stage of evaluation or tabulation of the data. Further, "sampling" is understood to refer to probability sampling only. Purposive sampling, and other systematic limitations of coverage are not considered^{1/}.

The reason for these various restrictions in the scope of this paper is to limit it to manageable length by concentrating on the most commonly met applications of sampling in Africa.

In any practical application of sampling the following topics have to be considered:

- Sampling units and sampling frames
- Multi-stage sampling and clustering
- Fixed versus variable sampling probabilities
- Stratification
- Sample size
- Sampling in time
- Relation to other sampling operations

^{1/} One further minor limitation of scope: no attempt is made to deal separately with problems raised by persons living in institutions, camps, hotels, etc.

These topics will be discussed in Section 2 independently of any particular application. They will then be taken up again in later sections in relation to their application to each of the three categories of sample enquiry listed above.

ABBREVIATIONS AND DEFINITIONS

- P.S.U. Primary sampling unit, i.e. first-stage unit.
- S.S.U. Secondary sampling unit, i.e. second-stage unit.
- I.U. Investigating unit. The smallest independent unit of personnel collecting data in the field. Typically consists of a single enumerator, but where enumerators work in teams whose members are interchangeable the IU will be the team - including the supervisor provided he works with one and only one team.

Unit of enquiry. Unit in respect of which information is collected. This may differ for different enquiries in the same survey.

Cluster sample. Sample in which the last stage of "sampling" is exhaustive, i.e. all the units of enquiry are included in the survey in each of the units selected at the next stage up. This usage is sanctioned by the Unesco supported Dictionary of Statistical Terms^{1/}, although the term is sometimes used to cover any multi-stage sample. Generally the term is used only where the clusters are geographical units. For example, a sample of households is not generally described as a cluster sample although it is in fact a cluster sample of individuals.

Self-weighting sample. Sample in which every unit of enquiry has the same overall probability of being selected. This means that data processing can be carried out as for a census, without differential weighting for different parts of the sample. Raising can be performed by multiplying the sample results

^{1/} Kendall, M.G., Buchland, W.R. (1957). Dictionary of statistical Terms. London/Edinburgh, Oliver and Boyd.

by a raising factor which is constant over the whole sample. Estimates which are ratios or percentages can be taken direct from the sample, without weighting.

P.P.S. sampling. Sampling with probability proportional to size. For example, villages may be selected with probability proportional to their census population. This will give larger villages a proportionally larger representation in the sample. The resulting bias in the sample is corrected either by introducing an inverse bias at a subsequent stage of sampling, or by re-weighting during data-processing.

2. GENERAL PRINCIPLES

2.1 Sampling units and sampling frames

It is well known that a sampling frame, or list of sampling units, is required to be accurate, exhaustive, non-repetitive (each unit appearing once only) and up-to-date. However, two of the most important requirements in practical applications in Africa are not explicitly mentioned in this standard formulation: (1) The units must be clearly and unambiguously demarcated - for example, for area units the boundaries must be clearly given, for social units such as households the membership criterion must be clear. (2) The units in the list must be traceable in the field: this means that, for example, there must not be so many villages of the same name, or so many heads of households of the same name, that it is impossible to identify on the ground the one selected on paper. Similarly, if houses are being selected in towns from an address list, the house numbers must not have been allocated so chaotically that the selected house cannot be traced. The use of aerial photographs for sampling of buildings is in most cases ruled out by this consideration; it is simply not possible for the average enumerator to identify on the ground a particular building indicated on a photograph.

The three main types of unit available for sampling are area units, units of housing or property, and households. The availability and convenience of sampling frames for such units will be discussed in later sections.

Variability in the size of sampling units may be an important consideration. It will be seen later in this paper that the organization of field work is often easier if the area sampling units are approximately constant in the population they contain; even if they are not, it is generally a help if their populations are known, even approximately.

Most African demographic surveys are concerned primarily with the estimation of certain rates or percentages. For this purpose no sample "raising" is involved: the estimates are taken direct from the sample,

after weighting if necessary. However, where estimates of aggregates are required, such as total population, it will generally be preferable to raise the sample by using supplementary information such as a census, even if this is out of date. Thus the sample data are multiplied by the ratio:

$$\frac{\text{Census population in whole domain}}{\text{Census population in sample}}$$

(preferably separately in each stratum). This is ratio estimation. For this purpose we need to know, once again, the approximate population of each primary sampling unit. This method of estimation gives a reduced sampling error. Alternatively, the same reduction can be achieved by taking account of the supplementary information at the sampling stage - sampling with probability proportional to size, or PPS sampling.

If supplementary information is not available (for example, if we use area sampling units not listed in the census tabulations), and if estimates of aggregates are still required, then we can only "raise" by multiplying the sample results by the reciprocal of the sampling fraction. This will generally give a larger sampling error, part of the error being due to variation in the population of the area sampling units. Sampling error will therefore be smaller if these units can be made approximately equal in population.

To summarize: a sampling frame should be accurate, exhaustive, non-repetitive and up-to-date. The units should be clearly demarcated and traceable on the ground. Ideally, units should be of approximately equal population; in any case, their sizes should be known - even an approximate knowledge of sizes is likely to be of value.

2.2 Multi-stage sampling and clustering

In any interview survey covering an area of more than a few square miles, some kind of grouping or concentration of the sample is a practical necessity in order to reduce the time spent by the IU (Investigating Unit) on travelling and other non-productive activities, in proportion to the time spent on data collection. Such grouping may involve sampling in two or more stages.

The effect of such concentration of the sample is generally to reduce sampling efficiency, when this is calculated on the basis of sampling error per unit of enquiry in the sample. This is because variability among units of enquiry is generally less within small areas than between them. Or, expressed in a less technical way, it is because members of one village tend to be similar, so that adding more people to the sample in the same village is less informative than going to another village. From the point of view of sampling error, therefore, the sample should be as widely spread as possible^{1/}.

However, once cost is taken into account, the advantage swings the other way. There is clearly a cost advantage in clustering, or grouping, the sample so as to reduce non-productive activities. Thus, travel costs increase with the number of sample villages, and at each village time must be spent "setting up shop" - that is, interviewing the chief, arranging accommodation for the IU, explaining the survey to the people. It is clearly wasteful to undertake all this if only a small amount of information is to be collected at each location. On the basis of cost alone, then, the more the sample is concentrated into a number of small areas or clusters the cheaper the operations will be. To some extent these arguments apply even for urban surveys and for sampling operations in which the IUs are not mobile (for the latter there will generally be need for a mobile supervisor).

Thus the requirements of sampling efficiency and cost operate in opposite directions and a balance must be sought between them.

A reasonable criterion of good sample design is minimum sampling error for given cost. This reflects both the above opposing tendencies. Using it, one can calculate an optimum degree of sample concentration or grouping, yielding minimum sampling error for given cost. (If we work in terms of

^{1/} It may be helpful to stress here that the problem involves both the concentration of the sample and the concentration of people's characteristics. The fact that people's characteristics are grouped in real life (neighbours tending to be alike) makes it desirable not to group or concentrate the sample, for the reason explained in the text above.

minimum cost for given sampling error, we reach the same optimum). If the two tendencies mentioned can be quantified the problem can be given a mathematical form and the optimum can be computed. The procedure is described in textbooks on sampling and is summarized below.

Four steps are involved:

- (i) Some measure has to be chosen representing the degree of grouping or concentration of the sample.
- (ii) This measure then has to be related mathematically to the sampling error.
- (iii) It also has to be related mathematically to the cost.
- (iv) The value of the measure then has to be computed which will minimize the sampling error for given cost.

These steps will now be considered in turn. We shall assume a 2-stage sample, or a 1-stage cluster sample^{1/}.

- (i) We shall measure the degree of clustering, or concentration, of the sample by n , the number of SSUs selected in each selected PSU.
- (ii) Assuming that the PSUs are large and equal in size and that the overall sampling fraction is small, it can be shown that the sampling variance is approximately proportional to

$$\frac{1}{mn} \left[1 + \delta(n-1) \right] \quad (1)$$

where m is the number of PSUs in the sample and δ is the so-called intra-class correlation. The latter is given by:

$$\delta = \frac{1}{N-1} \left[\frac{N\sigma_b^2}{\sigma^2} - 1 \right] \quad (2)$$

$$= 1 - \frac{N}{N-1} \frac{\sigma_w^2}{\sigma^2} \quad (3)$$

^{1/} The latter may be regarded as a 2-stage sample in which the 2nd stage sampling fraction is equal to 1. If a further stage of area sampling were included, making a 3-stage sample, there would be a small rise in sampling error and a small drop in costs. It is unlikely that there would be any important change in the optimum value of n , the number of ultimate stage sampling units selected in each PSU.

where N is the total number of SSUs in each PSU, σ_b^2 is the variance between PSU means, σ_w^2 the variance within PSUs, and $\sigma^2 = \sigma_b^2 + \sigma_w^2$ is the total variance. The intra-class^{1/} correlation δ reflects the extent to which the characteristic under investigation is concentrated or clustered, that is, how homogeneous it is within PSUs. This is apparent from equation (3) above: the more a characteristic is grouped according to PSUs, the smaller will be σ_w^2 , the within-PSU variance, and so the larger the value of δ . It is easily seen that if all members of any PSU have the same value of the characteristic, then concentration is extreme and $\delta = 1$. On the other hand if the characteristic is distributed at random without regard to PSUs, then $\delta = 0$.

For operational purposes in the present application only very rough estimates of δ are needed^{2/}. Actual values of δ met with in practice for different characteristics will be considered in a later section (3.2.1).

(iii) Cost function. It is probably fair to say that up to the present few African surveys have been limited primarily by cost. More often the bottleneck has been a combination of a time-limit and the availability of skilled field staff. For example, the survey has to be completed in time for the preparation of the national plan, or within a given financial year, or before the senior organizer's contract expires, or within the school holidays; and the work has to be supervised by senior field staff of whom only a limited number can be recruited at the level desired. Thus, in practice the limit within which the survey has to operate is a given number of IUs available for a given period. This leads to a "time" rather than a "cost" function, but the mathematics are the same.

1/ The term "class" may be misleading in the present application, where the class becomes the PSU.

2/ For a fuller discussion, and in particular for the computation of δ when PSUs vary in size, a recommended reference is: Hansen, M.H., Hurwitz, W.N., Madow, W.G. (1953). Sampling Survey Methods and Theory, (2 vols.). New York, Wiley.

The usual approach is to divide the cost, or available IU-hours, into two parts, one proportional to the number of PSUs selected, the other to the number of SSUs to be interviewed. The former will cover travel between PSUs and setting up shop, the latter the process of data collection. Let the time (number of IU-days) spent per PSU be T_1 and let the time per SSU be T_2 ^{1/}. The total available IU-days being constant, we have

$$T_1 m + T_2 mn = \text{constant}$$

or, dividing by T_2 and writing $T = T_1/T_2$,

$$T m + mn = \text{constant} \quad \text{-----} \quad (4)$$

T may be termed the "cost ratio" between first and second stage field work. Values met in practice in different circumstances will be discussed in a later section (3.2.2).

(iv) Minimizing (1) subject to the constraint (4) leads to

$$n_{\text{opt}} = \left[\frac{1-\delta}{\delta} T \right]^{\frac{1}{2}}$$

$$\div \sqrt{T/\delta}$$

It is easily shown that the optimum is very broad, so that substantial deviations from the optimal value of n do not affect the sampling error or cost very much.

Note that n is the number of SSUs selected in the sample in each PSU. If single-stage cluster sampling has already been decided upon, then automatically n is the cluster size, which is also the PSU size N . The problem of optimizing n is then the problem of the optimal size for PSUs.

^{1/} Strictly, these parameters should be understood as relating to the marginal time, that is, the time for each additional PSU or SSU included in the sample.

In fixing the sample allocation between the 1st and 2nd stage, the above cost and sampling error calculations are generally important and should always be attempted where the relevant data are available. However, there are many other considerations which may bear on the problem and which in some cases will be overriding. These will appear when we consider particular applications in later sections.

Once a value of n has been decided upon then, at least in principle, equation (4) determines the value of m , the number of PSUs to be selected, and hence the total sample size that can be covered with the available resources. This question is considered further in Section 2.5.

2.3 Fixed versus variable sampling probabilities

Fixed sampling probabilities make sampling easier and ensure simplicity at the data processing stage. However, variable sampling probabilities are often preferred for the following reasons:

- (i) To improve sampling efficiency by taking account of external information.
- (ii) To equalize the distribution of enumerator work-loads.

The first of these aims tends to be achieved by sampling PSUs with probability proportional to size (PPS). Other things being equal, the PSU total for any variable is likely to be approximately proportional to the size of the PSU (at least in some sense of the word "size"). When raising from a sample of PSUs, we normally divide each estimated PSU total by its corresponding sampling probability. If this probability is proportional to the PSU size, then the factor "size" appears in both numerator and denominator of the raised estimate and cancels. Thus variation in the estimates from different PSUs no longer reflects differences in PSU size. Such variation is essentially the variation which gives rise to sampling error. Thus sampling error is reduced by this procedure. However, this argument applies only where we are estimating aggregates. For estimation of rates or percentages there may be no improvement in sampling efficiency from PPS sampling.

Equalization of enumerator work-loads is achieved, assuming one enumerator works in each PSU, by sampling households within PSUs with probability inversely proportional to PSU size, since this procedure clearly leads to a constant sample of households in every PSU selected, provided "size" here is assumed to be measured by the number of households.

A popular sample design is to combine these two procedures. At the 1st stage, PSUs are selected with probability proportional to size. At the 2nd stage a fixed number of households is selected in each PSU. This gives both advantages simultaneously. Moreover, if the two measures of size used are the same, or at least proportional, then the 2nd stage sampling probability will be, in every PSU, the reciprocal of the 1st stage probability, so that the overall sampling probability (the product of the two) is constant. This means a self-weighting sample, which simplifies data-processing.

It is worth repeating that the above procedure aims for three advantages:

- (i) Improved sampling efficiency by taking account of PSU size (an important source of sampling error in the estimation of aggregates) at the sampling stage.
- (ii) Simplified field organization by ensuring a constant sample of households in every PSU selected.
- (iii) Simplified data-processing by use of a self-weighting sample.

In fact, however, these three advantages can only be secured simultaneously if the measure of "size" used for sampling at the 1st stage is proportional to the number of households currently living in the PSU. If no such measure of size is available, then one of the aims (ii) or (iii) cannot be exactly met.

It will be seen in later sections that we do not often have for every PSU, in advance of sampling, a measure of size which is sufficiently closely proportional to the number of households for the discrepancy to be ignored. We therefore have to choose between objectives (ii) and (iii).

It seems obvious that in the type of survey with which this paper is concerned the best policy will be to relax objective (ii). This is easily done by selecting at the 1st stage with probability proportional to size, using whatever measure of size is available, and at the 2nd stage sampling with the reciprocal of the 1st stage probability in each PSU. This ensures a self-weighting sample and gives an approximately constant sample of households in each PSU, with consequent convenience in field work.

The alternative policy of using fixed sampling probabilities at both stages abandons both objectives (i) and (ii). However, as we have seen in Section 2.1, if data on PSU size are used in ratio estimation, the equivalent advantages of (i) are secured at the data processing stage instead of at the sampling stage. Since the sample is also self-weighting, the only loss is then on objective (ii). This may not be serious, especially if PSUs do not vary too much in size. (Even if they do, it is often possible to reduce this variation by combining small units together and splitting large ones). Thus this very simple approach may often be preferable.

If a cluster design is used, with PSUs as clusters, then this means that the 2nd stage "sampling" probability is given as 1, which is fixed, so that the self-weighting design with variable probabilities mentioned above is automatically ruled out. The only choice is then between fixed and variable probabilities for the sampling of PSUs. In most cases it will be simpler to use a fixed probability, taking advantage of any available information on PSU sizes at the estimation stage only (ratio estimation).

Up to now we have assumed sampling in at most two stages. It may happen however that the smallest available area unit for which reliable boundaries can be identified is too big to serve as the field of activity of one IU. Rather than put several IUs into the same PSU, which would reduce sampling efficiency, it may be preferable to create secondary area sampling units in the field, in the sample PSUs only, of suitable size for one IU to work in. If a cluster sample of units of enquiry is desired, these created SSUs will be made of the size proposed for the clusters, so that the clusters will then be the SSUs. Otherwise, there will be a third stage of sampling, namely of units of enquiry within SSUs.

The same problems now arise again about the choice of fixed or variable sampling probabilities. If fixed probabilities are used throughout, this will be simpler at the sampling and data-processing stages, and supplementary information on PSU sizes can be introduced by ratio estimation with a consequent increase in efficiency of estimation of aggregates. On the other hand it is possible to achieve all three of the objectives listed above, as follows:

Select PSUs with probability proportional to the available measure of size. Create area SSUs within selected PSUs. Select SSUs with probability inversely proportional to the PSU size used at the 1st stage. Put one IU in each selected SSU. At the 3rd stage select units of enquiry with fixed probability (probability 1 if a cluster sample is desired).

This scheme uses knowledge of PSU sizes to increase sampling efficiency, it yields a constant work load for every IU provided the SSUs are created of constant size, and it is self-weighting. Of course in practice the SSUs will not be of exactly constant size, but the survey organizer is free to give whatever degree of emphasis he considers desirable to this requirement as the SSUs are specially created for the survey.

The practical procedure for selection of SSUs with probability inversely proportional to PSU size, as required by the above sample design, may be worth defining more exactly. Two methods are possible.

- (i) If the number of SSUs created in each PSU can be communicated to headquarters before sampling of SSUs begins, then the simplest procedure is that used normally for sampling with variable probability. Suppose PSUs are sampled with probability proportional to census population N_0 . Then the SSUs are required to be sampled with probability proportional to $1/N_0$. List the SSUs with the quantity $1/N_0$ against each. Cumulate these quantities. Select systematically (i.e. at a fixed interval, with a random start) in the cumulative column.

- (ii) If it is desired to perform the sampling locally in each PSU, then first fix a standard "expected" size for SSUs in terms of census population, say C . Then $N_2' = N_c/C$ is the expected number of SSUs to be created in a PSU, and N_2 is the actual number created. Round N_2' to the nearest integer, then select a random number between 1 and N_2' , another between $N_2'+1$ and $2N_2'$, another between $2N_2'+1$ and $3N_2'$, and so on. Any of these which fall between 1 and N_2 selects an SSU. This gives the desired selection probability, with an overall expectation of 1 selected SSU per PSU. This procedure is carried out for each PSU separately, with C constant for all PSUs. If a sample of k SSUs per PSU is desired, then read N_2'/k for N_2' in the above recipe when selecting the random numbers.

2.4 Stratification

The main purpose of stratification is to increase the efficiency of sampling. This purpose is achieved in two ways:

- (i) Stratification with unequal sampling fractions makes it possible to concentrate the sample towards areas where the variance is high or the cost low and hence to increase cost efficiency^{1/}.
- (ii) Stratification helps to ensure that the sample is well spread out over the domain of the survey. Stratification achieves this aim whether or not different sampling fractions are adopted for different strata.

For almost any area within Africa there is a wealth of geographical and anthropological studies published which give information relevant to the creation of strata. Unfortunately, very little of such information is statistical in form, and as a result there is little opportunity for making use of varying sampling fractions to improve the sampling efficiency of surveys, at least as regards the geographical stages of sampling.

^{1/} The same principle may be applied at any of the stages of a multi-stage sample, but in demographic applications it is likely to be limited to the first stage, and this is assumed in the discussion which follows.

Even if the desired information were available, there would be some doubt about how to use it, for two reasons: firstly, most surveys are multi-purpose, and the optimal sampling fractions for the different purposes are likely to differ; secondly, there is often doubt whether the objective is to minimize the error for the survey area as a whole or to produce acceptably low errors for individual domains of study within the survey area.

In view of these reservations, detailed discussion of this topic is perhaps unnecessary. However, we may mention the formula for optimization. The survey estimate as a whole has minimum sampling error per unit of cost when the sampling fractions for strata h are proportional to $\sqrt{V_h/C_h}$, where V_h is the variance and C_h the cost per unit of including a unit in stratum h . For a two-stage design, the relevant variance will be that between PSUs, and this will not generally be known, even for an attribute, with sufficient accuracy to distinguish between strata.

Turning to the second reason for stratification, the need to spread the sample, this is of considerable importance in areas as heterogeneous as are most African countries. In so far as the sampling frame is already in a stratified order, the same objective is achieved by systematic sampling (i.e. selection at a fixed interval in the list), the only disadvantage of the latter being that the extent of the reduction in sampling error which results cannot be satisfactorily estimated from the sample. With stratification, provided there are at least two sampling units selected in each stratum the sampling error can be estimated.

In general, there is a conflict between the desire to reduce sampling error by the maximum amount of stratification, and the desire to know the sampling error^{1/}. Two common solutions are the following:

^{1/} The practical importance of this problem is small because it only arises when stratification is pushed to the extreme. Several studies have shown that highly detailed stratification rarely pays off in terms of an appreciable reduction in sampling error. At the same time, when rigorous estimation of sampling error is impossible, approximate estimates are easily made (and commonly used) by ignoring part of the stratification process. Few statisticians would give priority to the rigorous computation of sampling errors.

- (i) Stratification is pushed as far as geographical knowledge allows, and sampling within the strata is then systematic. This gives maximum error reduction, with a compromise as regards knowledge of sampling error: we can estimate sampling error as far as the stratification takes us and this gives an upper limit, the further reduction in error due to systematic sampling being unknown.
- (ii) Strata, or zones, are created and exactly 2 PSUs are selected in each, by random selection. The sampling error can then be estimated satisfactorily, and with the minimum of difficulty since it is a simple function of the difference between the two PSUs of one stratum. In most demographic applications, with this sample design efficiency is maximal when the strata have equal populations. Further, the design is self-weighting only if the strata contain equal numbers of PSUs. The latter condition is likely to approximate the former and is perhaps the best basis for stratification (unless the self-weighting condition is met in some other way).

A further purpose which often leads to the introduction of stratification with unequal sampling fractions is the need to ensure an adequate representation of domains of study whose population is small. If a proportionate sample were selected, such a domain would be represented by a small sample, and the sampling error would be large for data relating to that domain^{1/}. If the survey objectives require reliable data for such domains they must be sampled with a higher sampling fraction and this means stratification with unequal sampling fractions.

Finally we mention the possibility of stratifying sampling units by size, with unequal sampling fractions in the strata. This is sometimes used as a substitute for PPS sampling. Instead of allowing the sampling probability to vary continuously in proportion to the PSU size, the PSUs

^{1/} Sampling variance is approximately inversely proportional to the sample size, not to the sampling fraction.

are grouped into size-classes and the sampling probability is made to vary step-wise, the probability in each stratum being proportional to the mean size for the stratum.

2.5 Sample size

Once the sample design has been fixed, including the number of sampling units at each stage to be selected within each unit of the next higher stage, the sampling standard error is approximately inversely proportional to the square root of the sample size, i.e. of the number of PSUs selected. In theory, the sample size could be decided by considering, in the light of the survey objectives, what would be an acceptable sampling standard error and fixing the sample size accordingly. However, the computation requires a knowledge of the variance in the population, more particularly the variance between PSUs. Generally this will be known only very roughly. The degree of precision which was adequate for fixing the allocation between sampling stages (Section 2.2) is hardly sufficient when it comes to fixing total sample size, which is inversely proportional to the variance itself.

However, as has been mentioned, most surveys in Africa are essentially limited by the size of the field force which can be adequately controlled: field supervision is the bottleneck. When the sample design has been fixed, it is therefore possible to compute the size of the survey which can be successfully fielded. Some very rough calculations should then be made along the lines indicated in the last paragraph to estimate the expected sampling error in order to see whether the survey is worth carrying out at all. This kind of approach to the problem of fixing sample size seems the most reasonable when knowledge of the relevant parameters is only very approximate; it is, in fact, how most African surveys have been planned.

When considering the acceptable limits of sampling error, it is of course essential to relate these to the amount of detail with which it is desired to break down the survey results. These two questions are opposite sides of the same coin.

2.6 Sampling in time

There would be little purpose in reporting results obtained from surveys unless the assumption could be made that such results would hold true over a longer period than that of the survey. In any field operation, therefore, an implicit assumption is made about the variates as a function of time. Moreover, any such operation itself constitutes some kind of sample in time.

It is as well to make these questions explicit and to ask, when designing a field operation, what are the assumptions regarding time variation and whether the sampling in time is adequate for the purpose envisaged.

Thus, if seasonal variations are believed to be important, the sample must systematically cover the seasons and in such a way that there is no confusion (interaction) between geographic changes, observed as a survey team moves from place to place, and true seasonal variation. Essentially the same applies where there are marked secular trends, even without seasonality. Further, if changes from year to year are important the survey must extend over a long period. It has been argued, for example, that demographic surveys designed to give a direct measure of vital rates have been of limited value in Africa because they have almost never covered more than one year and the rates vary widely from year to year. This type of objection should at least be considered before fixing the time schedule of any vital rate enquiry.

Another type of survey is designed to measure the effect on vital rates of some kind of official action or campaign. Here again, care is needed to eliminate seasonal effects. For example, if the survey is to be repeated at annual intervals then the field work must be performed at the same period each year.

2.7 Relation to other sampling operations

It is a common experience in all regions that, once an authority has decided to conduct a sampling operation, other interested parties try to join the bandwagon and to have additional sampling enquiries linked to the original project.

In Africa, demographic and housing surveys are particularly likely to find themselves associated with other types of survey. This is perhaps mainly because, in a region where sampling frames are rare, many surveys inevitably incorporate a round of field operations whose purpose is to list units in order to prepare a sampling frame. Such listing operations provide a convenient opportunity for collecting demographic or housing data.

In deciding whether to amalgamate operations which have different purposes, many factors have to be considered, of which the most important is probably the danger of overloading the learning capacity of enumerators. In this paper we limit discussion to those which directly concern sampling.

Firstly, we have to consider whether the units of enquiry for the different operations are compatible. If they are not, this is not necessarily an insuperable difficulty: it may be possible to cluster the units so that at least the clusters are compatible.

Secondly, are the desired sample designs compatible? Again, if they are not identical it may be possible to absorb one as a subsample of the other. If stratification requirements are incompatible it should be asked how much efficiency would be lost to either enquiry by deviating from the optimal stratification. Many demographic characteristics are distributed almost at random, so that it is often reasonable to subordinate the stratification requirements of a demographic enquiry to those of an associated agricultural or economic study.

Thirdly, are the time-sampling requirements compatible? For example, it may be decided that the most efficient plan for a demographic survey would be to cover the country progressively region by region, but for an agricultural survey all regions may have to be covered at harvest time. In many such cases the optimum solution is likely to involve separate rounds of field work for the different objectives.

3. SAMPLING OF RURAL SEDENTARY POPULATIONS FOR DEMOGRAPHIC CHARACTERISTICS

3.1 Sampling units and sampling frames

3.1.1 Area sampling units

In French-speaking areas of Africa satisfactory area sampling units are provided by the villages listed by the administrative census (or the full population census in the case of North Africa). These are usually defined with adequate precision for sampling. They are rather homogeneous in population size - averaging 300-350 population. The population given in the administrative census returns is in most cases somewhat inaccurate but can be used for PPS sampling or ratio estimation. However, if the administrative census population is used for PPS sampling at the 1st stage, and 2nd stage sampling is conducted on the basis of the fixed number of households selected in each PSU, it would not be advisable to assume self-weighting. An enquiry in Dahomey showed that the administrative census population was not in sufficiently close proportion to the number of households. (Such discrepancies may arise from error in the census, movement of the population since the census date, or variation in household size in different areas).

In English-speaking areas of West Africa, census enumeration areas (EAs) can be used in the same way. They average 1 000 population in Ghana and 300 in Nigeria. However, in the latter there are in many cases uncertainties about both the boundaries and the population of EAs and it would seem impossible to obtain reliable total estimates (as opposed to rates) for any variable.

In the countries of the East African Community, there are administrative units of average population 2 000. These are somewhat larger than is usually desirable for sampling but they have the advantage of reasonably well defined boundaries in most areas. In this part of Africa the rural population does not in general live grouped in villages.

In Lesotho, Malawi and Swaziland, EAs are very well mapped, with average populations of 1 000 for the first two countries and 500 for Swaziland. In Botswana the only available units are villages: these vary widely in size and, owing to seasonal migration, their population is known only on an approximate de jure basis. In Zambia, EAs of average population around 500 are expected to be created for the population census planned for early 1969.

In Ethiopia there exist units, variously named in different regions, whose identity seems to be fairly firmly established, although there are no maps. Their population averages perhaps a few hundred, but in most areas no useful estimates are available of the population of individual units.

In Sudan the smallest units are the "Sheikships". These are well defined administratively but not always geographically. Their populations in the 1955/56 census varied from about 100 to 1,000. Though unpublished, the figures are available at the Department of Statistics, Khartoum. In some areas tax lists might yield more up-to-date estimates.

In UAR the smallest area units are villages, with an average population of 2 - 3 000. These are very clearly delimited and their populations quite accurately known. It is possible that smaller EAs will be created for a census in 1970.

In many African countries (UAR is a clear exception), traditional units such as villages which ostensibly relate to an area may be regarded otherwise by the population. Thus, in the eyes of its inhabitants a village may be an extended family, acknowledging the authority of one chief. Anyone subject to this authority belongs to the village, wherever he may happen to live, and some members may live distant from the main area, surrounded by members of another village. Thus, villages frequently overlap. The attempt to treat such units strictly as areas, with distinct boundaries, may cause difficulties in the field. It is hardly possible to lay down rules for dealing with such problems. The statistician does his best in the local circumstances and tries to minimize the loss or rigour. In most cases he has to make some concession to popular opinion.

3.1.2 Sampling of housing and households

Except in UAR, no suitable sampling frames of housing or households exist in rural Africa. In principle the population census could be used for such purposes, but in practice it is hardly ever usable in rural areas as a sampling frame, whether for housing or for households. For housing, the units listed may be uncertainly defined, or they may be based on occupancy so that they change in precise coverage as people move in and out^{1/}, or they may be difficult to trace or identify. For a household sampling frame, the main difficulty is the instability of households over time, whether due to deaths, marriages or migration. This applies similarly to any attempt to use tax lists or administrative census records as household sampling frames.

Thus, in nearly all rural African surveys in which a sample of households is desired it has been found necessary to draw up new housing or household sampling frames by means of a special listing operation. The difficulties, mentioned above, which arise when a census is used as a sampling frame are then largely eliminated, mainly because the interval between listing and the subsequent survey can be made very short. In addition, the procedure is under the control of the survey organizer, who can arrange for any convenient definition of the listing unit; moreover, the identification of units selected can be assisted by the use of stickers affixed to houses and in many cases by employing the same unenumerators for listing and for the survey.

In only a few cases have African rural surveys been based on a census as a household sampling frame, and this has been only where the time interval between the operations was very brief.

^{1/} This may occur in error, contrary to instructions supplied to enumerators.

3.2 Multi-stage sampling and clustering

3.2.1 Variance between and within PSUs

As explained in Section 2.2, the degree of concentration or grouping of characteristics in the population - that is, the degree of correlation between members of the same PSU - is measured by the intra-class correlation, δ .

Values of δ have been investigated for a few characteristics in African surveys and censuses.

In Cameroun, birth and death rates were examined for 23 geographical strata in four regional demographic surveys. For birth rates, most strata showed values around .001 to .002, relating to area units of 350-400 population. In Ghana, calculations based on the number of children under 1 year old found in each census EA gave values of δ somewhat higher than the above although the area units concerned were larger: this may be due to recording errors in the census or to the fact that the number of surviving babies under 1 year old reflects infant mortality as well as fertility. It is probably safe to regard the Cameroun values as typical.

For death rates in Cameroun, somewhat higher values were found than for birth rates. δ averaged .002 but seldom exceeded .003.

Percentage employed among adult males was also investigated in the Ghana census data and very high values of δ were found, the average being about 0.1.

While few observational data are available, some general inferences can be made about the behaviour of δ in different circumstances.

- (i) In general, the larger the area unit considered, the smaller the value of δ to be expected, although the decrease in δ will be less than the increase in the size of the area unit. In practice, unless a very wide range of sizes is to be considered one may reasonably assume δ to be constant.

- (ii) Stratification at the 1st stage normally reduces the variance between PSUs and therefore reduces δ . The extent of the reduction depends on how closely the stratifying variable is related to the variable under study.
- (iii) Values of δ will naturally vary according to the variable studied, because some characteristics are more clustered than others. Age, sex and fertility may be expected to have low intra-class correlations, while for mortality the value should be higher (as was observed in Cameroun). Employment characteristics, tribe, nationality, religion, and other such variables in which people tend to sort themselves out into homogeneous groups, may show much higher values of δ unless the sample is stratified to reduce variation within strata for these particular characteristics.

Using these principles, results such as those quoted above, giving values of δ for a particular variable in a collection of surveys or strata, are precise enough to give a useful indication of the optimum sample size within PSUs, which is not very sensitive to the exact value of δ , once we have a suitable cost function.

3.2.2 Cost parameters

In Section 2.2 it was shown that optimal allocation of the sample between the 1st and 2nd stages may be computed by balancing a sampling error function against a cost function. In fixing the latter it was suggested that "cost" might best be measured in terms of IU-days^{1/}.

Two parameters then require to be specified:

T_1 = number of IU-days spent on each sample PSU, independently of the sample selected within the PSU

T_2 = number of IU-days spent per sample SSU over and above T_1 .

^{1/} IU: Investigating Unit. Either the enumerator, or the team of enumerators if they work together in teams in each PSU.

In Section 2.2 the parameter to be determined was n , the number of SSUs in the sample in each PSU, the SSUs being by implication households. In practical applications it is more convenient to measure n in terms of individuals than households. This is merely a question of multiplying n by the average size of households; if we measure the T s in terms of individuals also, there is no need for any change in the formulas quoted in Section 2.2.

T_1 has been found in African surveys to be typically 2 or 3 days. T_2 varies more widely; in particular, if the IU consists of a team of r enumerators T_2 will be r times as small. In theory one might expect an enumerator to be able to cover at least 30 individuals per day ($T_2 = 1/30r$), but in practice this level is not normally achieved over a whole survey and figures as low as 10 per day ($T_2 = 1/10r$) seem to be more normal. The team size r will generally be either 1 or around 5.

The formula for optimal n depends on T , the cost ratio T_1/T_2 . This will range roughly from 20 to 60 for enumerators working singly, while for enumerators working in teams of 5 the range is from about 80 to 300.

Taking $\delta = .002$ as a typical value for the intra-class correlation (see Section 3.2.1), this leads to the following optimal values of n :

Single enumerators: $n_{opt} = 100 - 170$

Teams of 5 enumerators: $n_{opt} = 200 - 400$

For age, sex and fertility, optima will be rather larger. For tribe, religion and employment they will be considerably smaller. Thus for a typical multi-purpose demographic survey it would be reasonable to select a sample of 1 or 2 hundred in each PSU.

For varying values of n , the cost for given sampling error is proportional to

$$\frac{1}{n} (T_1 + T_2 n) (1 + \delta n - \delta)$$

Example 1

Suppose $T_1 = 3$, $T_2 = 1/10$ (single enumerators), $\delta = .002$.

This leads to $n_{opt} = 122$.

If, instead of the optimum, we selected 5 times the optimal n in each PSU, then the cost to achieve the same sampling error would be increased by 51%. However, we can select twice the optimal n with an additional cost for given error of only 8%.

Example 2

Suppose $T_1 = 2$, $T_2 = 1/50$ (teams of 5 enumerators), $\delta = .0015$.

This gives $n_{opt} = 258$.

If we selected 4 times the optimal n in each PSU, then the cost for the same sampling error would be increased by 45%. If we selected twice the optimum, the cost for the same error would be increased by 10%.

Very broadly summarized, these results suggest that, for a typical rural demographic survey persons in Africa, it would be wasteful to select a sample of more than 300 or 400 per sample PSU if enumerators are to work singly, or about twice this number if enumerators are to work in teams of 5.

3.2.3 Cluster sampling

If the optimum sample size within PSUs is close to the total population of each PSU, then this is an argument for surveying everyone in each PSU, i.e. for cluster sampling.

A cluster sample may also be used even if the optimal n is much smaller than the PSU size, by creating secondary area units within PSUs equal in size to the optimal n . All households within each selected SSU are then surveyed, giving a cluster sample satisfying the optimality requirements.

There are a number of arguments which favour cluster sampling in demographic surveys.

- (i) If there is to be any sampling of households, a special listing operation is required in each selected area in order to prepare a household sampling frame, whereas if cluster sampling is used, no such operation is necessary. Against this, however, it may be argued that in many areas houses or households are clearly defined and where this is so the "listing" may be carried out concurrently with the survey: the enumerator merely selects for interview every n th house^{1/} or household. Other statisticians have argued that, even if there is to be no sampling of households, a house-listing operation ought to be regarded as an indispensable preliminary to any enumeration of the population. Either of these arguments tends to annul the argument at the beginning of this paragraph.
- (ii) Whatever is thought of the above arguments, it is certain that sampling requires additional instructions and training for enumerators and additional supervision to check error or falsification. In the case of concurrent sampling, unless the supervisor is actually present during the enumeration it may be impossible to detect the enumerator who deliberately numbers houses in such a way that the sample falls on smaller households. This error would have a direct distorting effect on the main survey results.
- (iii) The "household" is, in many African societies, difficult to define and rapidly varying in membership over a period of time. Thus, there is an advantage in avoiding its use as a sampling unit. Where the alternative of sampling houses or compounds is not satisfactory, cluster sampling provides the best solution. In any case cluster sampling is simpler for enumerators: by basing the enumeration on an area unit there is less doubt about who is supposed to be covered.
- (iv) In practice, it will be found embarrassing to apply any arrangement involving a sampling fraction within villages which is high

^{1/} The definition of a "house" for this purpose would depend on local circumstances.

(say over 0.5) but less than 1. The people not selected want to know why. Where the optimum sampling fraction within villages is above 0.5, this argues in favour of cluster sampling.

- (v) Efficiency in data collection may also argue in favour of a cluster design, particularly for vital rate surveys. Thus, if all households in a village are surveyed, there is more chance of getting information about births occurring in households which are absent from the village at the time of field work but which belong to the target population. There is also more chance of eliminating double counting of deaths due to reporting of the same death by more than one household.
- (vi) Finally, as we have noted above, substantial deviations from the optimum sample size per PSU can be accepted with little increase in cost. This fact may lend weight to the arguments quoted above for cluster sampling.

As long as the units available for use as PSUs are not much bigger than the optimum cluster, then these arguments are certainly overriding in favour of cluster sampling with PSUs as clusters. If, however, the PSUs are much beyond the optimum size, so that they are not suitable as clusters, then there is a choice between two policies: creating secondary area units for use as clusters, or sampling of houses or compounds. In most cases the latter solution will be both simpler to organize and more efficient as regards sampling error; the main objection is the difficulty of preventing manipulation by the enumerator. This solution has been adopted in several surveys in French-speaking African countries, but limited to large PSUs. In Cameroun, for example, in those PSUs whose population exceeded twice, thrice, etc., the optimal n , every 2nd, 3rd, etc. house was selected. In these surveys enumerators worked in teams under the direct supervision of a team leader (contrôleur). In these circumstances, perhaps only point (v) above remains as a substantial argument against this procedure. But where close supervision is not available it hardly seems advisable to adopt the house-sampling method as a regular procedure, in view of the danger of manipulation by the enumerator.

The alternative, creation and sampling of secondary area units, can be more easily controlled, although only at considerable additional expense. Two procedures are available which seem to be reasonably proof against falsification.

- (i) Secondary area units are created in selected PSUs by a special team in advance of the survey. A member of the team should accompany the survey enumerator when the latter starts work in the PSU in order to show him the area to be covered.
- (ii) SSUs are created by the survey enumerator, who prepares sketch-maps and descriptions. The supervisor then visits the PSU, undertakes the sampling, and checks over the boundaries of any selected SSU with the enumerator.

Both methods are clearly more expensive than the house-sampling procedure, the former because every PSU has to be visited twice, the latter because it seems impossible to avoid some wasted time while enumerators await the visit of the supervisor. A team-working method would of course solve this problem, since the supervisor could be present in the PSU throughout - but in this case the house-sampling procedure would be preferable.

3.3 Details of sample design

In the light of the principles outlined in Section 2.3 and the observational data reported in Sections 3.1 and 3.2, we can now suggest optimal sample designs for rural demographic surveys^{1/}.

The first step is to adopt a suitable area unit as PSU. This will generally be chosen as the smallest area unit for which reliable boundaries can be identified.

The next step is to estimate δ and T , from experience of other surveys conducted in comparable conditions, and hence to compute $n_{opt} = \sqrt{T/\delta}$, the optimum sample size per PSU. In assessing δ , a compromise will be necessary

^{1/} These designs ignore geographical stratification, which is treated in Section 3.4. If such stratification is used, they may be regarded as describing the design within a stratum.

since any real life survey has several objectives, for which δ will vary. Further, for estimation of T a decision has to be taken at this stage whether enumerators are to work singly or in teams. If singly, we may expect a value of n_{opt} around 100 -- 200; if in teams, about twice this.

Finally, a decision must be taken on the sampling procedure to be adopted within PSUs which are too large to serve as clusters: whether to create smaller area units to serve as clusters or to sample houses. The arguments for and against each method have been set out in Section 3.2.3. Broadly, the house-sampling method is simpler and cheaper but requires very close supervision. It would appear preferable where the team-working method of field organization is used and a supervisor is present in the PSU with the team at all times.

The detailed sample design can now be fixed. The basic principle will be to create, in each PSU, either clusters or house-samples of population very roughly equal to n_{opt} -- or at least not much more than twice this size. We now consider the various cases separately^{1/}.

3.3.1 Large PSUs to be split into smaller area units

We suppose here that a decision has been taken to deal with large PSUs by splitting them into secondary area units (SSUs) which will serve as clusters.

Splitting a PSU involves additional work. It is therefore reasonable to leave PSUs as they are unless their population exceeds about 2, or even perhaps 3, times n_{opt} . However, for PSUs which are so large that splitting becomes unavoidable, the number of SSUs to be created in the PSU should be based on the principle of making SSUs of size n_{opt} -- or perhaps a little larger, since the formula for n_{opt} ignores the increasing cost as the number of SSUs increases. In creating SSUs, the primary rule is that boundaries should be clear. Subject to this condition, it is desirable that the SSUs in a given PSU should be approximately equal in population.

^{1/} In all cases, unless otherwise stated, sample selection should be systematic rather than random.

We now further subdivide the discussion into two cases.

3.3.1.1 If PSU populations are known in advance of sampling

If no PSUs exceed about $3n_{opt}$, it would be better to abandon any subsampling and accept the PSU as the cluster in every case.

Otherwise, the larger PSUs should be split^{1/}. On the basis of the principles just stated, we can draw up a table showing the number of SSUs to be created according to the size of the PSU. The following example is based on the supposition $n_{opt} = 250$.

<u>Census population of PSU</u>	<u>No. of SSUs to be created in PSU</u>
< 600 (see note c)	1 (leave PSU unchanged)
600 - 750	2
751 - 1050	3
1051 - 1350	4
1351 - 1650	5
etc.	

- NOTES: a) Except at the beginning of the series, SSUs are created on the basis of 1 per 300 population, i.e. slightly more than n_{opt} .
- b) It may be advisable to allow for inaccuracy or obsolescence of the census data: n_{opt} should relate to the true current population.
- c) Very small PSUs (say, less than 100 pop.) may be grouped with neighbours if this is practicable, to reduce travel costs.

^{1/} The primary purpose of most demographic surveys in Africa is to estimate rates and percentages, not aggregates. Where, however, estimation of aggregates is important, it should be noted that the creation of secondary area units for which census data are not available will in general substantially increase sampling error, because it rules out the possibility of using the census directly for raising. In these circumstances, PSUs should probably be accepted as clusters up to a much larger limit than suggested above, or else the house-sampling method should be adopted, under proper control.

Now list the PSUs in the sampling frame, each with its census population, and, by referring to the above table, enter against each the number of SSUs to be created.

This amounts to a hypothetical list of SSUs, though they have not yet been created in the field. Sample SSUs from this list at a fixed interval from a random start.

Any SSU selected falls into a particular PSU. This gives a sample of PSUs.

Send enumerators into the selected PSUs to create the stated number of SSUs in each one.

Sample 1 SSU at random in each selected PSU.

Interview all households in each selected SSU.

The sample is self-weighting.

An alternative

The creation of a fixed and large number of SSUs in a given PSU is not a simple operation - e.g. to create exactly 10 SSUs in an area requires considerable skill by field workers. Thus, if the number of SSUs to be created per PSU is to average more than, say, 5 or 6 it may be advisable to adopt the following procedure.

Select PSUs with probability proportional to population. In each selected PSU create SSUs of population approximately n_{opt} (or a little more, as suggested above). Select among these with probability inversely proportional to the selection probability used at the 1st stage (for method, see end of Section 2.3). This gives a self-weighting sample and approximates to the method just described, but the number of SSUs selected per PSU will no longer be always exactly 1.

3.3.1.2 If PSU sizes are not known, even approximately

In this case we have no choice but to select PSUs with equal probability.

In each selected PSU, create SSUs of population approximately n_{opt} .

Select 1 SSU at random in each PSU.

Interview all households in each selected SSU.

The sample requires re-weighting at the data-processing stage.

An alternative

If there is to be a time interval between creation of SSUs in the whole stratum and the start of interviewing in that stratum, a simple alternative which yields a self-weighting sample is to collect the list of SSUs at a central point and sample from it with a fixed probability.

3.3.2 Large PSUs to be treated by sampling houses or compounds:

The principles here are much the same as in 3.3.1, except that, as long as satisfactory control of field work can be assumed, there is now no penalty in the sampling procedure as such, so that we can afford to introduce house-sampling as soon as the PSU population exceeds $1.5 n_{opt}$, and to aim at samples as close as possible to n_{opt} individuals in each PSU. We now consider two cases separately.

3.3.2.1 If PSU sizes are known approximately in advance of sampling

The simplest procedure would seem to be to select PSUs with PPS, then to sample houses (or compounds) within each PSU in proportion to the reciprocal of the probability used at the 1st stage, the constant of proportionality at the 2nd stage being chosen so that this procedure is expected to yield a sample of size n_{opt} persons in each PSU. The survey should cover all households in the selected sample of houses.

3.3.2.2. If PSU sizes are not known, even approximately

In this case we have to select PSUs with equal probability.

A rough reconnaissance will have to be made of each selected PSU in order to fix the 2nd stage sampling interval, which should be chosen so as to yield a sample of houses in the PSU containing approximately n_{opt} persons. This decision must be taken by a responsible officer. (The interval will, of course, vary in different PSUs).

The sample requires re-weighting at the data-processing stage.

3.4 Stratification

The main principles of geographical stratification have been outlined in Section 2.4 and there is little to add when it comes to practical application. Generally very little is known in advance about variances and costs in different strata, so that there is little scope for increasing sampling efficiency by sampling with unequal sampling fractions.

Except where the purpose of a demographic survey is to estimate the whole population of the country, there may often be a requirement for a sample design which gives equal precision in each of the main regions of the country. This implies approximately equal sample size in each region, which will usually mean unequal sampling fractions. Generally the most suitable plan will be to introduce these unequal fractions at the 1st stage of sampling.

Stratification and systematic selection can then be introduced as outlined in Section 2.4. If there is any sampling of houses or households, these should be selected by systematic sampling.

3.5 Sample size

When the principles outlined in Section 2.5 are applied in practice they have in most cases led to a sample of 50 000 to 150 000 persons for African rural demographic surveys. Most such surveys have, however, been regional in coverage. If it is desired to cover a whole country whose total population exceeds 5 million, such a sample would probably not give sufficient detail for small regions.

It is hardly possible to go further than this and to specify more precisely what sample size is necessary for a satisfactory survey. As explained in Section 2.5, in most cases the procedure has been to carry out the largest survey which can be financed and supervised. The analysis of the data in terms of small regions is then carried as far as sampling error permits.

3.6 Sampling in time^{1/}

In any demographic survey there are at least three different aspects of the time sampling which require a decision at the planning stage.

3.6.1 Concentration or dispersion of the period of field work

Some African demographic surveys have been organized in such a way as to complete the field work (at least for any given round) in the shortest possible time. Examples are the Nigerian demographic survey of 1965/66 and the Ghana population survey of 1966. This procedure involves the use of a large force of enumerators, with supervision by personnel who have a relatively low level of training. The alternative, widely used in French-speaking African countries, is to use a small force of enumerators divided into a few teams, each closely supervised by a well trained team leader (contrôleur), which move around the country collecting data over a period of many months.

The former procedure produces data which are more convenient for the demographer, since the whole area of enquiry is covered simultaneously although there is the disadvantage that seasonal trends are not allowed for; in the latter method, difficulties of interpretation can arise through confusion between seasonal movements of the population and movements of the interviewing teams. On the other hand, the precision of the data collected by the second method is likely to be superior, because supervision is stricter and enumerators gain efficiency through experience; in addition, the smaller number of enumerators may make it possible to apply stricter standards in recruitment. In one regional survey in Cameroun, for example, almost every questionnaire was inspected within 24 hours of its completion by one of the team leaders who moved with each of the ten survey teams. Any detectable errors were rectified at once by a return visit to the household. Every questionnaire was also

^{1/} The subject matter of this section has been treated in fuller detail in the ECA Seminar on Vital Statistics held in Addis Ababa in December 1964. Among the seminar papers, see in particular: on recall lapse, Technical paper on non-sampling errors and biases in retrospective demographic enquiries (E/CN.14/CAS.4/VS/3); on follow-up surveys, Methods of obtaining vital data in developing countries (E/CN.14/CAS.4/VS/5); on use of long-term retrospective data, Uses of census or survey data for the estimation of vital rates (E/CN.14/CAS.4/VS/7).

inspected within a short interval by one of the two qualified statisticians who organized the field operations. In these circumstances enumerators learn rapidly from their mistakes. Supervision of this quality is not possible with the "one-shot" type of survey.

3.6.2 Retrospective versus follow-up survey

If data are collected on vital events, a decision has to be taken whether to obtain these by retrospective questioning, in which respondents are asked to recall events occurring within a specified period, or by a follow-up survey, in which changes in the population are directly observed by re-visiting a given sample of households after an interval. The former method is known to produce serious recall error, though techniques have been suggested for adjusting for this by assuming a particular mathematical model for recall lapse, and thence extrapolating to zero recall-period where recall error is assumed to vanish (field work must be spread over a year to eliminate seasonal variation)^{1/}. Inevitably there is uncertainty about the validity of this latter assumption and about the choice of model, but the method has the advantage of cheapness since only one round of field work is needed. The follow-up method involves fewer assumptions but costs about twice as much at the field stage and in most cases involves much greater complexity at the processing stage^{2/}. Moreover, it cannot altogether avoid the need for retrospective questioning since there is no other way of obtaining information on babies who are born and die between visits (a check on pregnancies at the 1st round may, however, eliminate some errors of this kind). In a typical follow-up survey the first and last rounds are separated by 12 months; there may or may not be an intervening round between these two. The method was used very successfully in Nigeria in 1965/66.

1/ Som, R.K. (1968). Recall Lapse in Demographic Enquiries. Asia Publishing House, Bombay.

2/ The Nigerian survey of 1965/66, however, showed that it is possible to arrange the work in such a way that data-processing is no more complex than for a single round survey.

3.6.3 Long-term retrospective data

If the right questions are asked it may be possible to get reliable data covering a much longer retrospective period than one year. Thus, women may be asked to state the total number of children ever born to them. Brass^{1/} has suggested relatively elaborate methods of analyzing such total fertility data, in conjunction with data relating to the last 12 months, in such a way as to eliminate both reporting error and bias due to differential survival of mothers. The number of assumptions made, however, is considerable. Mortality may similarly be estimated by asking respondents whether their parents are living.

Methods of this kind have the advantage of covering a longer period and so eliminating short-term fluctuations in natality and mortality - which may be particularly important in countries of low rainfall. Their disadvantage is that they involve considerable manipulation of the data, depending on numerous assumptions; often conclusions are considered to be supported by the convergence of evidence from more than one type of analysis. It is very difficult to know exactly how much confidence can be justifiably accorded to evidence of such a complex character.

It should be noted that collection of long-term retrospective data is in no way an alternative to collection of short-term information. In most surveys both methods have been used.

3.7 Relation to other sampling operations

This topic has been adequately discussed in Section 2.7. When a proposal is under consideration for linking a demographic survey to another sampling operation, the implications should be carefully studied, both from the sampling and non-sampling point of view, and a decision taken on the merits of the case. The factors to be considered as regards sampling have been outlined in Section 2.7.

^{1/} Brass, W., et. al. (1967). The Demography of Tropical Africa, Princeton.

4. SAMPLING OF URBAN POPULATIONS FOR DEMOGRAPHIC AND HOUSING CHARACTERISTICS

4.1 Sampling units and sampling frames

4.1.1 Area units

Most large African towns have been covered by a complete census during the last 10 years. Such an operation commonly leaves behind a legacy of a well defined structure of enumeration areas (EAs) which are generally suitable for sampling purposes, at least in the larger cities. In the smaller towns (say, under 100 000 population), such units may not be sufficiently numerous for satisfactory sampling, but in many such cases the small size of the town makes a complete count feasible even in the context of a national demographic sample survey, so that no sampling problem arises.

Many African towns also have available recent large-scale aerial photographs or detailed maps which enable blocks to be marked out which can be used as area sampling units. In general, such units have the disadvantage that their populations are not known even approximately; but in some cases rough population estimates can be obtained by counting buildings observed in the photograph or map and multiplying by an estimate of the ratio population/buildings obtained from the latest census for each district of the town. Such estimates can be used either directly for sample raising (ratio estimation) or for creating blocks of approximately constant size.

4.1.2 Units based on housing and property

In most African towns there are numerous conceptual difficulties in defining units based on housing, and this means that it will generally be troublesome to try to sample from lists of such units obtained during a census^{1/}. Sampling buildings from aerial photographs is even less practicable.

^{1/} However, the method seems to have been used with success in UAR.

A few African towns, however, are laid out in a simple grid structure which leads directly to clearly defined spatial units of property, or "lots". Brazzaville and Kinshasa are notable examples, where the lots are termed parcelles. Such units may be sampled from aerial photographs, from census records (Brazzaville) or from municipal registers (Kinshasa). In some towns this may be possible in certain districts only (Khartoum, Omdurman).

Finally, in many African towns even where lay-out is unsystematic, municipal records exist which are claimed to cover every inhabited building. These should never be used for sampling without a preliminary field check to confirm (a) that a very high proportion of the inhabited units are in fact listed and (b) that any selected unit can be unambiguously and rapidly identified on the ground.

4.1.3 Household lists, tax lists, etc.

African urban populations generally have a very high level of mobility (see Section 4.6) and household lists are not likely to be useful for sampling when they are more than a few months old. This means that in practice household sampling frames are virtually unavailable.

In a number of countries tax lists are available which are supposed to cover the whole urban population. This claim should be checked by a small sample field operation before such lists are used for sampling. Even if the list is found to be complete, there may be difficulty in converting it from a list of taxpayers into one of households. Recently in Tananarive, for example, such a list was examined carefully with a view to its use as a household sampling frame; over a period of 2 weeks, new categories of exempted persons constantly came to light and the difficulties of converting to a household list were seen to be more and more troublesome. Finally the attempt to use the tax list was abandoned as impracticable.

4.2 Clustering and multi-stage sampling

There seems to be almost no evidence on variances between small areas and field costs for the urban sector in Africa^{1/}.

Urban populations, of course, tend to stratify themselves by income and social class and such strata are readily identifiable. Within any stratum it is reasonable to suppose that demographic characteristics would be rather evenly distributed, i.e. the intra-class correlation would be low. This would favour a large sample in each area unit, or large clusters, but on the other hand the cost advantage of grouping the sample is clearly very small in urban areas, so that it is not immediately obvious whether the optimum degree of grouping would be greater or less than in the rural sector. However, it is doubtful whether this approach - balancing a continuous cost function against a variance function - is really appropriate in urban conditions because there is practically no cost advantage in creating household groupings or clusters larger than can be covered by one enumerator in one day. This is because enumerators must go home for the night and it is just as easy for them to go to a new location in the morning as to return to the old one. Thus, a reasonable policy is to fix the number n of households to be selected per PSU at less than the daily enumerator's quota. This will almost certainly be less than the optimal n computed without taking account of the fact that the enumerator returns home each night, so that such a policy will be optimal as regards sampling error and cost.

The above discussion assumes that a suitable sampling frame exists for sampling households, or at least some kind of unit based on housing. As we have seen in Section 4.1.2, this situation is comparatively rare. At this point it is necessary to complete the discussion under two separate heads.

^{1/} An exception is the Ghana census data on percentage employed, mentioned in Section 3.2.1, where the intra-class correlation between EAs was found to be at least as high as for rural areas - namely over 0.1.

4.2.1 Cities or districts of regular lay-out

In areas organized on a grid-plan it is a simple matter to create blocks of almost any desired size. A possible design would be (4.2.1.1) to define blocks corresponding approximately to the size which can be covered by one enumerator in one day and to use these as clusters, covering every lot within the selected block. This would meet the optimality requirement already noted. If, however, the overall sampling fraction is fairly high, the cost advantage of such grouping of the sample is quite small and the question arises whether it would not be preferable to select (4.2.1.2) a single stage sample of lots, with no grouping or clustering.

In choosing between these two alternative plans the main factors to be considered are probably the following:

- (i) Can lots selected from the sampling frame be readily traced on the ground? If much time is likely to be wasted on this activity a cluster sampling procedure with blocks as clusters is likely to be preferable.
- (ii) Is an important objective of the survey to collect data on housing? Informal observation suggests that housing characteristics are highly grouped or clustered and there will be an advantage in dispersing the sample as much as possible. This argues in favour of a single stage sample where housing data are important to the survey. The same reasoning applies where employment data are considered important.

4.2.2 Cities or districts of irregular lay-out

If lay-out is irregular, then even if an accurate house-list exists it is unlikely to be convenient for sampling, for two reasons: (i) the size of the houses (or whatever unit is listed) is likely to be excessively variable, and (ii) selected houses are likely to be difficult to find. However, provided neither of these conditions holds the case may be treated as similar to that in the preceding section.

If there is no convenient house-list, at least four alternative plans may be considered.

4.2.2.1 Aerial photographs or maps may be used to create small blocks which will serve as clusters. If the blocks are too large the sampling error is likely to be excessive, but if they are small the amount of work involved in block creation for a fair-sized town may be considerable. This method was used in Yaoundé in 1964 (population 100 000), where blocks of 40-50 houses were created. From the sampling point of view this seems likely to have been well above the optimal block size, but the listing of houses and households within selected blocks served at the same time to provide a sampling frame for a household budget survey in which one enumerator was to operate in each block. The block size was determined primarily with this purpose in view. Moreover, the creation of blocks of substantially smaller size would have required a prohibitive amount of labour. Incidentally, housing information was collected in this survey from the household budget sample, not at the house-listing stage.

4.2.2.2 An exhaustive housing census, or a complete count of buildings or houses, may be carried out which will provide a sampling frame for a demographic enquiry. This procedure was used in Addis Ababa in 1967. The difficulty of identifying selected units may be reduced by the use of stickers which are affixed at the initial listing stage. The method would not give a very satisfactory sample in towns such as those of southern Ghana and Nigeria, where buildings vary greatly in size, unless it is found possible to list relatively small housing units rather than buildings (but this has its own difficulties - see Section 4.2.2.4 below).

4.2.2.3 In larger towns it may be convenient to introduce a 1st stage of area sampling by selecting a sample of EAs. There is then a choice between dividing the selected EAs

into blocks by a field operation, each selected block to serve as a cluster, or listing all houses in selected EAs and sampling from this list. The latter is considered in Section 4.2.2.4. In the former case, blocks should be made as nearly constant in population as practicable, without sacrificing the requirement of well defined boundaries.

4.2.2.4 Creation of blocks is not a simple operation and it is probably no more laborious to list every house in the selected EAs and sample from this list. There are two alternatives regarding the choice of listing unit. We may choose a large unit, such as the compound or house number, or we may choose a smaller unit intended to correspond approximately to one household. In some areas the compound typically corresponds to one or two households, and in such cases this is almost certainly a better solution than 4.2.2.3. Where there are no clear-cut compounds, however, there will be difficulties with either type of unit. The large unit is likely to be excessively variable in size^{1/}, while any small unit will (i) be difficult to re-identify for the enumerator who is assigned to survey the selected sample units and (ii) may cut across households, so that demographic data are not easily obtained for the selected sample. Whether these difficulties are over-riding can only be determined by preliminary trial and error in the areas concerned. If they are, cluster sampling by creation of blocks becomes unavoidable. Cluster sampling is also strongly indicated for the case of a follow-up survey covering a period of a year or more: household composition, and utilization of

^{1/} In urban areas of English-speaking West Africa, buildings (or "house numbers") are common in which over 100 people live, while in the same areas one also finds many traditional single-family houses.

housing structures, is so fluid in many urban areas that clear delimitation of sampling units over a period of time can only be achieved on an area basis.

4.2.3 Note on housing units

The African Recommendations for the 1970 Housing Censuses^{1/} define a housing unit as "a separate and independent place of abode intended for habitation by one household, or one not intended for habitation but occupied as living quarters by a household". It is explicitly provided that one housing unit may be occupied in fact by more than one household, and one household may occupy more than one housing unit. It follows that housing units in this sense cannot serve as sampling units for a survey which covers both housing and demography unless special provisions are made for re-weighting. It also seems likely from this definition that in many cases the "housing unit" will not be very clearly demarcated and there may be differences of interpretation between the units listed by one enumerator and the selected sample units subsequently surveyed by another. In most cases, then, it will be desirable to avoid use of the "housing unit" as a sampling unit.

Where sampling of housing is still desired one may seek a larger unit, more clearly defined and not cutting across households or housing units. As stated in earlier sections, this may be the lot or the compound in cities where these are clearly defined, or elsewhere it may be the house number. The latter has the disadvantage, in many cities, of being highly variable in size besides being difficult to trace.

There is of course no objection in principle to use of different concepts as sampling units and reporting units. The United Nations definitions are concerned with the latter.

^{1/} E/CN.14/CAS.5/CPH/10 (July 1967). Economic Commission for Africa.

4.2.4 Note on housing surveys

Housing characteristics are generally highly clustered, so that a dispersed sample is desirable. In some cases it may be sufficient to collect housing information from only a subsample of the units counted during a listing operation. In some surveys, housing data have been collected from a subsample of the demographic sample. This procedure involves omitting unoccupied housing, though it is possible to make separate arrangements to cover this. It may also lead to difficulties over incompatibility of housing and demographic units.

4.2.5 Summary

It will be seen from previous sections that we can choose from many satisfactory designs for demographic and housing surveys in urban areas. The choice depends on objectives and requires a careful study of the situation in the town concerned. It is, of course, possible that a different choice would be made for different districts of the same town.

The designs that have been suggested above are the following:

..... Exhaustive census

4.2.1.1 Blocks created on basis of grid-plan of city. Cluster sample, using blocks as clusters.

4.2.1.2 Single stage sample of lots selected from grid-plan of city.

4.2.2.1 Blocks created from aerial photographs. Cluster sample, using blocks as clusters.

4.2.2.2 Exhaustive housing count. Sample of compounds or houses from this list.

4.2.2.3 Sample of census EAs. Creation of blocks within selected EAs. Cluster sample, using blocks as clusters.

4.2.2.4 Sample of census EAs. Listing of houses or compounds within selected EAs. Sample from this list.

4.3 Fixed versus variable sampling fractions

Of the sample designs just listed, only the last two require a decision as regards variability of sampling fractions; the others either involve no sampling or it is self evident that sampling should be with fixed probabilities (at least within strata).

For designs 4.2.2.3 and 4.2.2.4 there is a choice between sampling EAs with fixed probability or with probability proportional to their census population. If fixed probabilities are used at the first stage it would be reasonable to use fixed probabilities also at the second stage in order to simplify data processing. EA population figures from the census would then be used for ratio estimation if estimates of aggregates are desired. If, on the other hand, EAs are selected with probability proportional to census population, it will be convenient to sample at the second stage with the reciprocal of this probability. In case 4.2.2.3 the technique used would be one of the two methods described at the end of Section 2.3. In case 4.2.2.4 it will be simpler just to compute the second stage sampling fraction in each EA (as the reciprocal of that used at the first stage, multiplied by a constant for all EAs) and select the houses or compounds accordingly.

The PPS design leads to more complex sampling but simpler processing. It also tends to equalize the work loads in each EA. It seems marginally preferable to the alternative fixed probability scheme.

4.4 Stratification

Demographic characteristics can be expected to vary considerably between different districts of any city (although not, perhaps, between neighbouring small areas). For housing characteristics, variation will be even greater. Such natural strata are readily identifiable and it will in most cases be easy to stratify a town by rough income level and social class on a geographical basis.

As the rich are, on any normal classification, less numerous than the poor, it may in certain cases (depending on survey objectives) be desirable to increase the sampling fraction for wealthier areas in order to ensure

that they are represented by an adequate sample. In a housing survey, other special groups may require an augmented sampling fraction for the same reason.

Stratification may also be used for spreading the sample more evenly over the domain of study, but systematic selection achieves this aim more easily (see Section 2.4). If there is to be any sampling of households or units of housing, these should normally be selected by systematic sampling.

4.5 Sample size

Demographic sample surveys in African urban areas have generally been based on a sampling fraction of about 1/10. Except for very large towns this gives too small a sample for reliable data on rare events (births and deaths), or on very clustered characteristics (such as employment data) if multi-stage sampling is used. However, any increase of the sampling fraction above about 1/5 would raise the question whether an exhaustive census would not be more fruitful. For this reason most urban demographic surveys have been carried out either in cities of over half a million population (Kinshasa, Addis Ababa), or as part of a national demographic survey, in which case data on individual towns are not given (Ghana, Nigeria, Morocco). In a few additional cases, the demographic interest has been subsidiary, demographic data being collected in a listing round designed primarily to provide a sampling frame for a household survey (Yaoundé, Libreville, Freetown). But in most medium and small sized African towns demographic data have been collected only by exhaustive census.

4.6 Sampling in time

Any study of the demography of modern African towns must take account of the very high mobility of such populations. Some idea of the significance of this movement is given by the following table based on follow-up surveys conducted in Abidjan (1963) and Yaoundé (1964/65). An additional complication found in Abidjan was the existence of a marked seasonal variation in migration.

Status at time of follow-up	Per cent of initial population	
	Abidjan 1 year later	Yaoundé 6 months later
Still in the same dwelling	74 %	84 %
In another dwelling in same town	17 %	12 %
Left town	7 %	4 %
Dead	2 %	0.5 %
	<hr/> 100 %	<hr/> 100 %
New arrivals	18 %	7 %
Births surviving	5 %	2 %
	<hr/> 23 %	<hr/> 9 %
Net increase	23-7-2=14 %	9-4-0.5=4 %

Source: Based on Tables 35 and 36 in Démographie Comparée (1967), 7 : Déplacements temporaires et Migrations. I.N.S.E.E., Paris.

One result of this high mobility is a large gap between the de facto and de jure populations. Moreover the exact de jure population will vary substantially depending on the residence qualification.

All this means that careful attention has to be paid to the definition of the target population in terms of residence. The objectives of the enquiry should first be carefully considered. Then, either definitions should be adopted appropriate to these objectives, or data should be collected for the widest possible group (i.e. the de facto population) together with information on length of residence, so that a classification by different residence periods can be tabulated.

In view of the seasonal variability of migration, a survey spread over a whole year is likely to lead to some inconsistencies unless care is taken to distribute the sample widely over the town at all times (i.e. a systematic movement of the enumerators over the town, as a group, should not be allowed).

Follow-up surveys in urban areas have been attempted in Abidjan and Yaoundé. They can give valuable data on mobility, as we have seen in the above table. However, they require very careful attention at the data-processing and analysis stage in order to separate the different categories of movement of population and to draw valid conclusions. Cluster sampling is desirable in order to define unambiguously the sample to be covered.

Long-term retrospective data (total fertility) have been collected in urban areas as in rural (c.f. Section 3.6.3). However, correction for reporting error and differential migration does not appear to have been attempted for urban populations and the very high rate of migration would appear to make this difficult.

4.7 Relation to other sampling operations

As has been mentioned in Section 4.5, many urban demographic surveys in Africa have been linked to other operations, either larger or smaller.

On the one hand, they may constitute part of a national demographic survey. In this case it is usual to create an urban stratum in which the sample design may be quite different from that used in the rural sector. The main relevance of the national operation for the design of the urban operation is that in most cases where a national sample survey is taking place the urban survey is required to give data for the urban sector as a whole, and not for individual towns. This, of course, affects the sample size (see Section 4.5). The existence of the national survey is relevant also to timing: if the urban and rural operations follow a very different time schedule there is likely to be difficulty over the interpretation of data on rural/urban migration.

On the other hand, urban demographic surveys are sometimes carried out almost incidentally, as the first round of a household survey. In this case it is to be expected that the needs of the latter will largely determine the sample design.

Housing surveys are rarely carried out as independent operations. Usually they are combined with a demographic survey or census, sometimes with a household survey. The main difficulties arising from these combinations have been mentioned above in Sections 4.2.3 and 4.2.4.

5. CIVIL REGISTRATION ON A SAMPLE BASIS FOR RURAL POPULATIONS

5.1 General

The setting up of a civil registration system on a representative sample basis, designed ultimately to expand to total coverage, has been much discussed but no African country has yet instituted such a scheme^{1/}. This section deals with the sampling requirements for a project of this kind. The discussion does not follow the form of previous sections since the problems are rather different.

If registration is carried out by visiting all households in a sample ("active registration"), we have an operation which is, from the sampling point of view, no different from the demographic surveys discussed in Section 3. In the present section we shall therefore deal only with the case of an essentially passive registrar, who waits for the report of a vital event to reach him, at most visiting hospitals, clinics, etc., but not households.

We have introduced the terms "active" and "passive" registration here to distinguish the registrar who is required to visit regularly every household in his area from the registrar who is not so required. Naturally, the latter will not in fact be inactive, but the term "passive" appropriately reflects the fact that in this case the ultimate responsibility for the reporting of vital events lies with the public. The distinction is important from the sampling point of view because the "active" registrar can, during his visits to households, count the base population which he is covering. For the passive system, the estimation of the base population presents a special statistical problem which distinguishes the case from the vital rate sample survey covered in previous sections of this paper.

^{1/} See in particular E.C.A. African Seminar on Vital Statistics, held in Addis Ababa in December 1964. (Report sales No. 65.XVII.6)

This limitation of the discussion in this section is in no way intended to imply that the approach to vital rate estimation through permanent surveys, or active registration techniques, is likely to be any less fruitful. Such a policy has been followed in India (annual series of single round surveys) and for a period in Pakistan and Thailand (follow-up surveys supplemented by vital registration). In Africa, there have been moves in this direction in Senegal and UAR^{1/}. However, any civil registration scheme which it is intended ultimately to extend to 100 per cent coverage on a continuing basis could hardly be based on active registration.

The discussion below is also limited to the rural sector, because it seems likely that, in most case where a civil registration sample scheme might be set up, the urban areas would be set aside for total, rather than sample, coverage.

5.2 Type of sample required

When a sample civil registration scheme is set up, with the intention that it shall ultimately be extended to total coverage, it may be supposed that legal arrangements will be made to cover the whole country, at least nominally, from the start. The sampling relates not to the legal framework but to the operational efforts which will be made to get the system working efficiently. These will be specially directed towards a sample of registration areas.

For practical and administrative reasons, the registration areas are likely to be fairly large administrative areas, each covering a population of several tens of thousands. All residents of the sample areas will be covered by the scheme, so that we have a cluster sample. Population data for computation of vital rates will come from an independent source (see Section 5.3) so that we may reasonably adopt equal probability sampling

1/ World Population Conference, 1965, Vol. III, United Nations. Cantrelle, P.: Repeated demographic observation in a rural area in Senegal, p. 200. Vukovich, G.: The UAR project for measuring vital rates in rural areas, p. 195.

for the clusters. Systematic sampling from a list arranged geographically will ensure that the sample is well dispersed over the country. The total sample size in terms of persons covered will presumably run into several hundreds of thousands.

5.3 Base population

The only serious statistical problem is how to obtain the base population, that is, the population of the sample area.

There is very little evidence about the reliability and durability of census figures in Africa. It is known that for small areas (EAs) the population figures are often very inaccurate or very quickly go out of date but they seem likely to be more reliable for larger areas. Moreover, much of the error is probably random so that for a large sample the net error may be small. We need an accuracy of perhaps $\pm 5\%$ (corresponding to ± 2.5 per thousand on a birth-rate of 50 per thousand)^{1/}. Whether this can be obtained by extrapolating from the last census, using other relevant information where available, is a question which has to be answered for each individual country in the light of all available knowledge.

If it cannot, a sample survey will have to be envisaged within the registration sample. We now consider the sample design for such a survey.

To obtain a good estimate of total population, the census figures will have to be used for ratio estimation, and this probably means that the survey PSUs must be census EAs. The size of the sample required for the survey will then depend on the correlation between census population for each EA and the population found in the same EA during the survey. In some cases analysis of an earlier survey might supply an estimate of this correlation. If this is not possible, some kind of sequential sampling procedure would seem the best policy: a guess is first made as to the sample size required, field work is carried out on this basis, results are analyzed and if the sampling error is found excessive the sample is

^{1/} The intervals quoted are intended as corresponding to a high confidence level -- 90 or 95 per cent.

enlarged and further field work is done (this second round might follow one year after the first). Sampling of EAs for each round should be systematic. (Care will be needed in dovetailing the two systematic samples). The sample will, of course, always stay within the areas selected for the registration sample.

It may be noted that the carrying out of a sample survey at the same time as the introduction of sample registration serves also another purpose, namely a check on the accuracy of registration. The survey can collect independent data on births and deaths and thereby provide an independent estimate of the vital rates. Case-by-case matching of vital events found by the two operations may also be possible, at least on a subsample. For this checking purpose it is clearly desirable that the survey should be based on a cluster sample. The clusters will then be EAs. The fact that these clusters are smaller than the registration areas will cause some difficulty with matching. Locality of residence, and if possible the complete address, should be recorded in the registration, but in some cases the EA to which this corresponds may not be clear, and it may be necessary to limit the matching purposely to exclude EAs where such difficulties are likely.

Finally, comparison of the survey population estimate with the preceding census should give an improved idea of how to extrapolate the population estimate into future years - an essential exercise if the sample registration area is to continue to give useful estimates of vital rates.