

**PROPERTY OF THE  
PUBLICATIONS BRANCH  
EDITORIAL LIBRARY**

# **Sample Design and Estimation Procedures For a National Health Examination Survey of Children**

A description of the sample design used for the second cycle of the Health Examination Survey.

DHEW Publication No. (HSM) 72-1005

---

**U.S. DEPARTMENT OF HEALTH, EDUCATION, AND WELFARE**  
Public Health Service

Health Services and Mental Health Administration  
National Center for Health Statistics  
Rockville, Md. August 1971



# NATIONAL CENTER FOR HEALTH STATISTICS

THEODORE D. WOOLSEY, *Director*

PHILIP S. LAWRENCE, Sc.D., *Associate Director*

OSWALD K. SAGEN, Ph.D., *Assistant Director for Health Statistics Development*

WALT R. SIMMONS, M.A., *Assistant Director for Research and Scientific Development*

JAMES E. KELLY, D.D.S., *Dental Advisor*

EDWARD E. MINTY, *Executive Advisor*

ALICE HAYWOOD, *Information Officer*

## OFFICE OF STATISTICAL METHODS

MONROE G. SIRKEN, Ph.D., *Director*

E. EARL BRYANT, M.A., *Deputy Director*

## DIVISION OF HEALTH EXAMINATION STATISTICS

ARTHUR J. McDOWELL, *Director*

PAUL T. BRUYERE, M.D., *Deputy Director*

HENRY W. MILLER, *Chief, Operations and Quality Control Branch*

JEAN ROBERTS, *Chief, Medical Statistics Branch*

LINCOLN I. OLIVER, *Chief, Psychological Statistics Branch*

JAMES T. BAIRD, Jr., *Acting Chief, Methodological Research Staff*

Vital and Health Statistics-Series 2-No. 43

DHEW Publication No. (HSM) 72-1005

*Library of Congress Catalog Card Number 76-608549*

# PREFACE

This report presents a detailed description of the sample design and estimation procedures employed by the Health Examination Survey in a nationwide survey of youths 6-11 years of age in the non-institutional population of the United States. The objective of the survey was to collect data which would provide national estimates and distributions of various health characteristics related to the growth and development of this target population.

The overall responsibility for the development of the design and other sampling aspects of the survey was that of Walt R. Simmons, Assistant Director for Research and Scientific Development, National Center for Health Statistics (NCHS). Garrie J. Losee, formerly Assistant Statistical Advisor, NCHS, with assistance from George A. Schnack, Office of Statistical Methods, NCHS, shared in the planning of the design and was responsible for the development and execution of specific sampling procedures. Innovations in the design, such as the Latin-square modification of the controlled selection techniques, are the joint contribution of all three above-named persons. The Statistical Methods Division, Bureau of the Census, particularly Robert Hanson, devised the techniques for, and performed the ultimate stage selection of, sample segments from 1960 census listings.

This report was prepared jointly by the three staff members listed as its authors. Much of the report is based upon internal unpublished documents written by Messrs. Losee and Simmons.

# CONTENTS

	Page
Introduction -----	1
Preliminary Considerations and Specifications-----	2
Determination of Sample Size-----	3
First- Stage Design and Selection of PSU's-----	5
General -----	5
Formation of HES Superstrata-----	6
Selection of First-Stage Units-----	11
Within PSU Design-----	16
Problems of Development-----	16
Coverage of the Universe -----	16
Selection of Localities Within PSU's-----	17
Selection of Segments-----	19
Selection of Sample Children-----	21
Estimation Procedure-----	22
Variance Estimation-----	25
Background-----	25
Summary of Applicable Theory-----	25
Application to Cycle II Data-----	27
References -----	31
Appendix I. Glossary of Terms-----	33
Appendix II. Procedure for Forming and Stratifying PSU's in the Current Population Survey and the Health Interview Survey Designs-----	34
Formation of PSU's-----	34
Stratification of PSU's-----	34
Appendix III. Household Questionnaire-----	37

### SYMBOLS

Data not available-----	---
Category not applicable-----	...
Quantity zero-----	-
Quantity more than 0 but less than 0,05----	0.0
Figure does not meet standards of reliability or precision-----	*

# SAMPLE DESIGN AND ESTIMATION PROCEDURES FOR A NATIONAL HEALTH EXAMINATION SURVEY OF CHILDREN

E. Earl Bryant, *Office of Statistical Methods*, and James T. Baird, Jr., and Henry W. Miller, *Division of Health Examination Statistics*

## INTRODUCTION

The Health Examination Survey is one of the major survey programs employed by the National Center for Health Statistics to obtain information about the health status of the U.S. population. It is a part of the National Health Survey, authorized in 1956 by the 84th Congress as a continuing Public Health Service activity.

The National Health Survey employs three different survey programs to accomplish its objectives.<sup>1</sup> One of these is the Health Interview Survey in which persons are asked to give information related to their health or to the health of other household members. The second program, Health Resources, obtains health data and health resource and utilization information through surveys of hospitals, nursing homes, and other resident institutions and through the entire range of personnel in the health occupations. The third major program is the Health Examination Survey (HES).

The Health Examination Survey collects data from samples of the civilian, noninstitutional population of the United States and, by means of medical and dental examinations and various tests and measurements, undertakes to characterize the population under study. This is the most accurate way to obtain diagnostic data on the prevalence of certain medically defined illnesses. It is the only way to obtain information on unrecognized and undiagnosed conditions—in some cases, even nonsymptomatic conditions. It is also the only way presently available to obtain

distributions of the population by a variety of physical, physiological, and psychological measurements. Although the sample is designed primarily to estimate the prevalence of specified health and health-related conditions in the population, the design also makes possible the study of relationships of the examination findings to one another and to certain demographic and socioeconomic factors.

Successive and separate survey programs are conducted for specific age segments of the population. These programs, referred to as "cycles," are concerned with certain specified health aspects of that subpopulation. Thus, the first cycle of the Health Examination Survey was conducted between November 1959 and December 1962 and was directed toward the civilian, noninstitutional population from ages 18-79 years inclusive. The examination was focused primarily on certain chronic diseases, principally cardiovascular diseases, arthritis and rheumatism, and diabetes. Also included were a dental examination, tests for visual and auditory acuity, an X-ray, electrocardiographic tracings, blood chemistry tests, and numerous body measurements. The sample size of this cycle was 7,710 persons, of which 6,672 (86.5 percent) were examined. Details of the plan of this initial program<sup>2</sup> and reports of various methodological studies<sup>3-11</sup> and of the findings<sup>12</sup> relative to that cycle are available.

The target population of the second cycle of the Health Examination Survey consisted of children ages 6-11 years inclusive. This cycle

became operational in July 1963 and was concluded in December 1965. The primary focus of the examination was on various parameters of growth and development, but it also screened for heart disease; congenital abnormalities; ear, nose, and throat diseases; and neuro-musculo-skeletal abnormalities. The size of the sample of this cycle was 7,417, of which 7,119 (96.0 percent) were examined. Several methodological reports,<sup>13-17</sup> as well as reports of findings,<sup>18</sup> have been published, and others are being prepared.

A detailed report of the plan, operation, and response results of the second cycle has also been published.<sup>19</sup> While that report does include a general description of the sample design, it was necessarily limited by the scope of the report. It will, therefore, be the object of this report to describe in detail the various aspects of the sample design and selection procedures, weighting techniques used for population estimation, and procedures employed for variance estimation.

## PRELIMINARY CONSIDERATIONS AND SPECIFICATIONS

The development of a successful sample design must take into account all relevant factors and circumstances. In view of the primary mission of the Health Examination Survey, this means that there must be a blend of primary survey objectives, budgetary resources, logistical considerations, time limitations, organized speculation concerning population parameters, and unit operating costs. These and other requirements in Cycle I dictated that a highly stratified multistage probability type of design be used in contrast to some possible alternative of a more subjective or volunteer selection of examinees.

The similarity between Cycles I and II, particularly with respect to their broad mission, indicated a similar probability type of design for Cycle II. It should be pointed out, however, that while of necessity several features were common to both designs considerable statistical exploration was carried out for Cycle II to determine the optimum design with respect to sample size, sample allocation, sampling frame, and operational procedures.

In the early planning stages of Cycle II, two problems basic to the sample design received considerable attention. These were the age segment of the population to be examined and the sampling frame to be used. The original concept was that the age group to be studied in Cycle II would be persons ages 6-17 years inclusive. As the detailed planning proceeded, however, it became apparent that the differences between persons in different age segments of this population group were so great that separate programs were required. Therefore, it was decided to redefine the Cycle II target population as children from ages 6-11 years, inclusive, and to follow this program with a third cycle which would have youths 12-17 years, inclusive, as its target population.

Since almost all the population in the age group 6-11 years are in school for a large part of the time, it was felt that a sample design which used the school populations as an element of stratification might have some operational advantages. For example, if schools could be grouped by type (public, parochial, private, etc.), size, socioeconomic characteristics of the students enrolled, and segregation factors, a sample of children from one or more schools in each group might minimize the number of specific locations from which the sample children would come. Although some consideration was given to using the schools in this way as a sampling frame, the idea was abandoned because of the unavailability of the necessary classificatory data concerning the schools, difficulties anticipated during summer months, and geographic coverage of nonpublic school children.

Consideration was also given to selecting an original sample of 15,000 to 25,000 children and to making some of the simpler elements of the examination on all. A smaller sample would be selected from the original group and would be subjected to the additional examination and tests requiring more elaborate equipment or procedures. Important advantages of such a scheme were that it would permit a two-phase selection of the smaller sample and would provide poststratifying information that would reduce sampling variance. This plan was discarded, however, because of the operational problems it seemed to present.



In the final analysis, the sample design of Cycle II was developed essentially from a set of specifications which took into consideration requirements and limitations placed upon it. It was important that the requirements be consistent with survey objectives and that the limitations not be so serious as to materially distort the objectives. Specifications of primary importance were as follows:

1. The target population would be the non-institutional population of the United States from 6-11 years of age, inclusive, with one exception. Because of operational difficulties experienced in Cycle I, all children residing upon any of the reservation lands set aside for the use of American Indians would be excluded.
2. The data collection mechanism developed and proved during Cycle I would be used, with appropriate modifications. Examinations would be conducted in mobile examination centers, two of which would be in operation simultaneously in different parts of the country.
3. The total period of data collection for Cycle II would be between 2 and 3 years. Other time limitations were a maximum 6-day workweek, a 5-week-per-year loss of time due to vacations and holidays, and a 7-day loss per move from one examining location to another.
4. The length of an individual examination would be between 2 and 3 hours. Approximately 12 children would be examined per day.
5. Experienced and qualified personnel in the field staff for Cycle I would be retained to the extent necessary to perform the data collection operation in Cycle II.
6. The schedule of examining locations or stands must take into account the climate, especially to avoid conducting the survey in Northern States during the winter.
7. Certain cost factor limitations such as budget loads projected for each of the fiscal years 1962 and 1963 must be observed.
8. The examination objectives would be concerned primarily with factors of physical and mental growth and development.

9. Ancillary data would be collected through the use of questionnaires. These would consist of a household questionnaire, a medical history of the child completed by the parent, and an interviewer-administered medical history questionnaire. Also, a questionnaire would be sent to the school at which the sample child is a student.
10. Maximum target tolerances for sampling variability would be set for several key statistics, permitting a general analysis by broad geographic regions, population size groups, and other major subgroups such as age, sex, and limited socio-economic factors.

## DETERMINATION OF SAMPLE SIZE

The size of sample required for a survey is influenced by a number of factors. These include the sample design, estimating procedure, confidence-tolerance specifications, variability and prevalence of population characteristics to be measured, available budget and unit costs, and operational constraints placed on the design. Once all such factors are determined, and therefore fixed, the sample-size requirement for a stratified design will vary depending on how the sample is allocated to strata and how the sample is clustered within strata. In designing Cycle II, one such factor examined was how to allocate the sample in such a way as to produce estimates with minimum variance for a fixed budget.

One of the design specifications was to permit analysis by broad geographic regions. Thus, a first consideration was to divide the population of the United States into a number of geographic regions approximately equal in population size. As explained in greater detail in a later section, this resulted in four regions with further stratification occurring within each. The latter stratification further produced an equal number of strata within each region which in turn were also approximately equal in size. Under these conditions, population variances are often about the same magnitude in each stratum. Also, the cost of examining an individual should be somewhat similar from one examining location to another.

These features of the design indicated that an equal allocation of the sample strata would

be approximately optimum. Thus, in determining the sample size, the main consideration was how to allocate the sample between the first- and second-stage units, that is, the number of primary sampling areas or units (PSU's) and the number of sample persons per PSU.

To determine an optimum solution to this problem, a cost relationship,  $B = C_0 + C_1 m + C_2 m \bar{n}$  was assumed where  $m$  = number of PSU's,  $\bar{n}$  = number of sample persons per PSU;  $B$  = total budget for the survey,  $C_0$  = overhead costs,  $C_1$  = costs associated with a PSU such as travel between PSU's, and  $C_2$  = costs associated with persons such as cost to examine a person. The optimum values of  $m$  and  $\bar{n}$  for a two-stage cluster sample design which yield estimates with minimum variance for a fixed budget are:

$$\bar{n} \text{ (optimum)} = \frac{S_w}{S_b} \sqrt{C_1/C_2}$$

$$m \text{ (optimum)} = \frac{B - C_0}{C_1 + C_2 \bar{n} \text{ (optimum)}}$$

where  $S_w$  and  $S_b$  are components of the total population standard deviation due to variation within PSU's and between PSU's respectively. Estimates of  $S_w$  and  $S_b$  were computed from data collected in a probability sample of 14 PSU's completed early in Cycle I, using the formulas:

$$\hat{S}_b^2 = \frac{m}{\sum_{i=1}^m} \frac{(P'_i - P')^2}{m-1}$$

$$\hat{S}_w^2 = \frac{m}{\sum_{i=1}^m} \frac{n_i}{\bar{n}} P'_i Q'_i$$

where  $n_i$  is the actual number of sample persons in the  $i^{\text{th}}$  PSU.

The proportion of the population with a specified health characteristic,  $P'$ , components of variances  $\hat{S}_b^2$  and  $\hat{S}_w^2$ , sampling error of the estimated proportion  $\sigma_{P'}$ , and optimum values of  $m$  and  $\bar{n}$  are shown in table A for a number of health conditions. The information on which these estimates were based was not ideal for designing the Cycle II sample, since it related

to adults and not to children; and the health characteristics were not the same as those to be considered for Cycle II. Thus the assumption had to be made that the information about variances and unit costs for the survey of adults also held approximately true for growth and development characteristics of children 6-11 years of age.

As is nearly always true in surveys which have multiple objectives, the optimum values of  $m$  and  $\bar{n}$  vary for the different variables. For the statistics proposed for this survey, the values ranged from 57 PSU's and 105 sample persons per PSU for estimating diabetes to 95 PSU's and 35 sample persons per PSU for estimating peripheral vascular disease. The choice of a best design was not possible because all variables were of equal importance, and therefore a compromise had to be made. If precision and budget were the only factors to influence the choice of a best design, possibly the choice would be about 75 PSU's and 64 sample persons per PSU since the optimum for eight of the variables requires 75-95 PSU's and a similar number requires less than 75 PSU's. Sample designing is not that simple, however. The best design is also a function of things other than sampling error, such as availability of personnel and equipment and procedures which minimize measurement errors.

For the Health Examination Survey, an item of considerable importance and concern is non-response. It was learned in Cycle I that a high cooperation rate can be expected, however, if one is willing to make several callbacks to find the family at home and to set a time for the examination that is convenient for the sample person. To accomplish this requires that the examining team remain in the area at least 2 or 3 weeks. Another important factor which influenced the choice between design alternatives was the need to minimize the loss of effective time resulting from moving from one location to another. Thus there is a limit to the number of PSU's that can be completed with available resources and time limitations.

As seen in table A, for a 40 PSU design it is possible to examine 180 persons per PSU, or a total of 7,200 persons, for about the same

Table A. Comparison of a 40 PSU design with minimum variance, fixed-cost optimum designs for 14 health statistics collected in Cycle I of the Health Examination Survey

Health statistics in Cycle I	Proportion of population with characteristic ( $P'$ )	Within PSU variance ( $S_W^2$ )	Between PSU variance ( $S_B^2$ )	Optimum design			Selected design		
				Number of PSU's ( $m$ )	Number of persons per PSU ( $\bar{n}$ )	Sampling error $\sigma_p'$	Number of PSU's ( $m$ )	Number of persons per PSU ( $\bar{n}$ )	Sampling error $\sigma_p'$
High blood pressure-----	.168	.135	.00468	87	45	.0095	40	180	.0115
Organic heart disease---	.084	.076	.00147	77	60	.0060	40	180	.0070
Peripheral vascular disease-----	.105	.089	.00514	95	35	.0090	40	180	.0120
Arthritis-----	.215	.162	.00664	90	41	.0110	40	180	.0135
Visual acuity-----	.278	.198	.00297	72	69	.0090	40	180	.0100
Edentulous persons-----	.169	.136	.00439	86	46	.0090	40	180	.0115
Weight greater than average-----	.605	.235	.00411	75	64	.0100	40	180	.0115
Diabetes-----	.017	.017	.00011	57	105	.0015	40	180	.0020
Headaches-----	.743	.191	.00180	64	87	.0080	40	180	.0085
Nose bleeds-----	.113	.100	.00093	64	87	.0055	40	180	.0060
Tinnitus-----	.327	.220	.00248	67	79	.0090	40	180	.0095
Dizziness-----	.431	.245	.00658	83	51	.0115	40	180	.0140
Orthopnea-----	.076	.067	.00095	71	71	.0050	40	180	.0055
Chest pains-----	.310	.214	.00423	77	60	.0100	40	180	.0115

cost as that for the optimum designs indicated in the table. Although the sampling errors are larger for a 40 PSU design than for the corresponding optimum design, for most practical considerations in using the results of the survey the 40 PSU design and the optimum design can be viewed as having about the same reliability. Therefore, when all factors were considered, the 40 PSU design was chosen as best under prevailing circumstances.

## FIRST-STAGE DESIGN AND SELECTION OF PSU'S

### General

A major and often expensive task in designing and implementing a national population survey is to establish and maintain a sampling

frame containing the target population, to order the population in such a way that facilitates sample design efficiency, and to select the sample units. Fortunately, much of this work had already been done as part of the U.S. Bureau of the Census Current Population Survey (CPS) and the Health Interview Survey (HIS). For these purposes, the 3,103 counties and independent cities which compose the total land area of the United States had been combined into 1,891 primary sampling units and had been further stratified into 357 homogeneous classes or strata. The first-stage sample units for both CPS and HIS (at the time of designing Cycle II) contained 357 PSU's, one from each stratum.

To implement these surveys the Bureau of the Census maintains a trained field staff of several hundred people located in 12 regional offices. The Bureau also maintains a continuing

program for keeping the sampling frame current through the collection of building permits issued in the sample PSU's. Thus design efficiency was significantly enhanced by taking advantage of the Bureau resources in designing Cycle II.

The Cycle II sample of PSU's consists of 40 of the 357 HIS PSU's. It is not a subsample in the usual sense of the word, however. The characteristics of the 357 HIS sample PSU's were used as a matter of convenience, to collapse the 357 HIS strata into 40 HES superstrata. Then by use of controlled selection, one HIS stratum, referred to subsequently as "first-stage units" or FSU's, was selected from each superstratum with probability proportional to the size of the first-stage unit. Finally, the sample PSU that originally represented the HIS stratum was chosen for the Cycle II sample.

Although detailed descriptions of the HIS and CPS sample designs have been published,<sup>20,21</sup> a brief summary of how the PSU's were formed and stratified is presented in appendix II to facilitate understanding of the full design of Cycle II.

In this section the procedures for forming superstrata and for selecting first-stage units from superstrata are discussed.

### Formation of HES Superstrata

To understand how superstrata were formed it is useful to view all of the PSU's in an HIS stratum as a single unit. In this report, these units are called first-stage units since the first stage of sample selection in Cycle II was of FSU's. The first step in the Cycle II design was to stratify the 357 FSU's into 40 superstrata on the basis of the characteristics of the HIS sample PSU's. This was done in a manner which maximized the degree of homogeneity within superstrata with respect to FSU population size, geographic proximity, degree of industrialization, and degree of urbanization. Stratification was carried out within 16 mutually exclusive cells formed by classifying the FSU's into four population density classes within each of four geographic regions of the United States.

Other features of the design which had an influence on how the superstrata were to be formed

included the need to produce self-weighting estimates, to produce estimates for each of the four regions, and to have a sample of approximately the same size for each PSU. The implications of these conditions on design efficiency are that the regions should be about the same size, each region should contain about the same number of strata, and each stratum should contain about the same number of people. This type of balance was achieved by creating 10 superstrata in each region with the condition that each of the population density classes (largest standard metropolitan statistical areas (SMSA's), other large SMSA's, other SMSA's and highly urban counties, and rural and other urban areas) would also contain 10 superstrata.

To create regions containing about the same number of people, it was necessary to redefine the commonly used Bureau of the Census regional boundaries. A comparison of the two definitions is shown in figure 1.

- The four geographic regions are:
- Northeastern-- identical to the Census-defined Northeast Region.
  - Midwestern--- Census-defined North Central Region less Kansas, Nebraska, North Dakota, and South Dakota.
  - Southern----- Census-defined South Region less Oklahoma and Texas.
  - Western----- Census-defined West Region plus those parts detached from the North Central and South Regions.

Figure 1 is somewhat misleading, however, in that the actual content of the Cycle II regions does not follow the State lines in all instances. This is the result of assigning FSU's to regions according to the State within which the sample PSU in the HIS design was located. Some strata in the Western Region contain PSU's actually located in the Midwestern and Southern Regions. Similarly, some strata in the Midwestern and Southern Regions include PSU's located in the Western Region. The problem is not serious, however, since only a very small proportion of a region's population is involved in the overlap.

The four population density classes, which also divide the country into four roughly equal parts, were defined on a sliding scale. For example, the Atlanta SMSA in the Southern Region

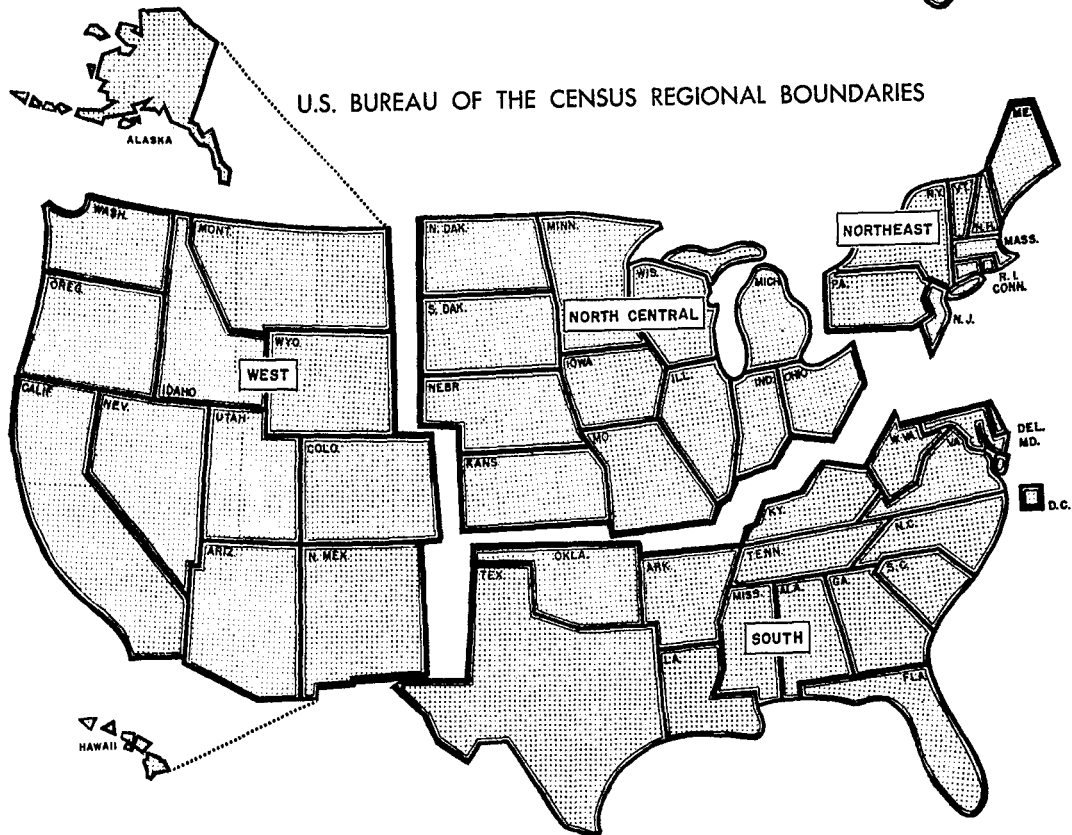
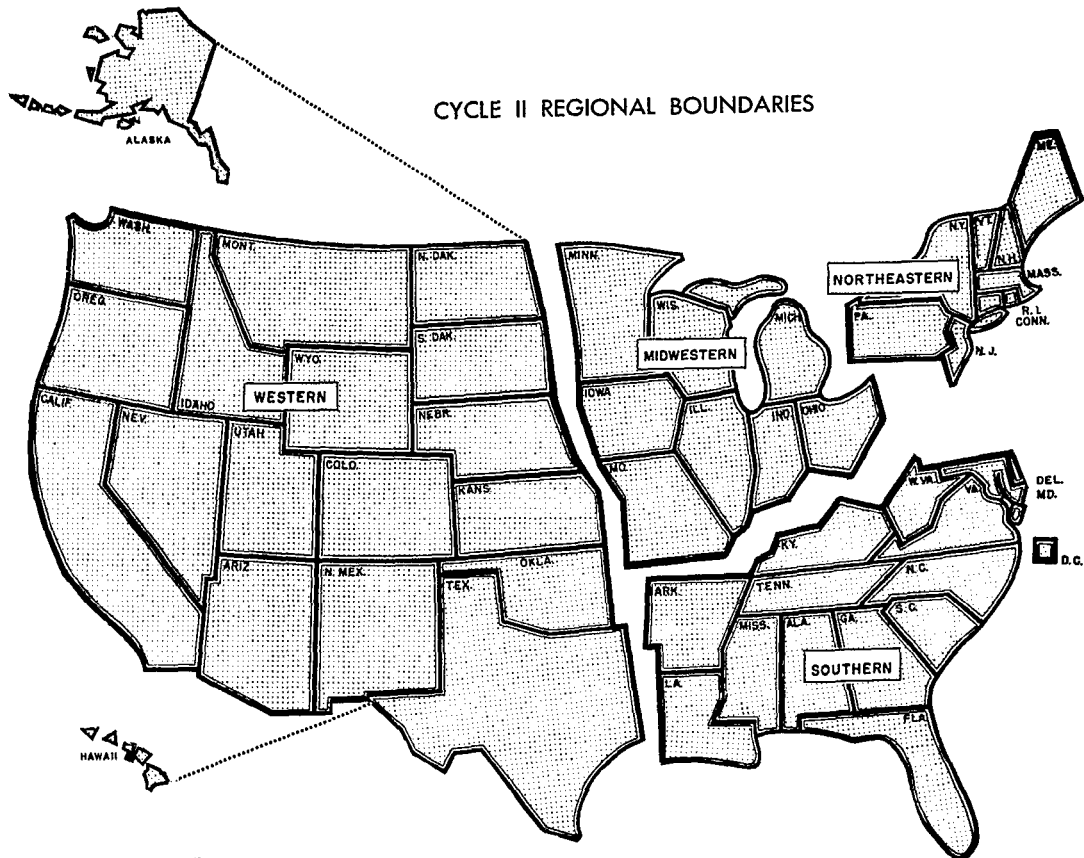


Figure 1. Comparison of Cycle II regional boundaries with those defined by the Bureau of the Census.

with a population of about a million people is equated on the scale to Philadelphia in the Northeastern Region and to Chicago in the Midwestern Region. The reasoning is that Atlanta has a position of economic importance in the Southern Region similar to that of the other two cities in their respective regions. The approximate population ranges for size classes are shown in table B.

The average size of superstrata and the distribution of FSU's and superstrata by geographic region and population density class are shown in table C.

Note that each density class within a region was represented by either two or three superstrata and that the average size of the superstrata was around 4.5 million people.

Seven of the superstrata were self-representing. That is, each contained a single FSU.

The New York SMSA was split to form two superstrata, as was Los Angeles. The others were Detroit, Philadelphia, and Chicago SMSA's.

The non-self-representing superstrata were formed by grouping two or more FSU's. To the extent possible the FSU's in a superstratum were similar in size, as well as in other characteristics mentioned above.

In the highest two population density classes, the FSU's tended also to be self-representing. In general these were SMSA's of more than 500,000 people, as indicated in table C. Superstrata composed of "other SMSA's and highly urban counties" in the Northeastern and Western Regions were, in the most part, made up of self-representing FSU's as shown in table C. Contrastingly, all FSU's in the Southern Region and 85 percent in the Midwestern Region were non-self-representing.

Table B. Definition of population density classes within geographic regions

Geographic region	Population density classes			
	Largest SMSA's	Other large SMSA's	Other SMSA's and highly urban counties	Rural and other urban areas
Northeastern-----	SMSA's with more than 3 million people	SMSA's with 1-2 million people	SMSA's with less than 1 million people	All rural and other urban areas
Midwestern-----	SMSA's with more than 3 million people	SMSA's with 500,000-2 million people	Other SMSA's and highly urban counties with less than 500,000 people	All rural and other urban areas
Southern-----	SMSA's with more than 700,000 people	Other SMSA's	Non-SMSA, highly urban areas	All rural and other urban areas
Western-----	SMSA's with more than 1,100,000 people	SMSA's with 500,000-1,100,000 people	Other SMSA's	All rural and other urban areas

Table C. Distribution and average size of superstrata and first-stage units by geographic region and population density class and whether or not self-representing

Geographic region and population density class	Number of superstrata			Number of FSU's			Average size of superstrata	Average size of FSU's
	Total	Self-representing	Non-self-representing	Total <sup>1</sup>	Self-representing	Non-self-representing		
							(In thousands)	
United States---	40	7	33	364	122	242	4,462	492
Largest SMSA's-----	10	7	3	16	16	0	4,419	2,762
Other large SMSA's----	10	0	10	64	54	10	4,269	667
Other SMSA's and highly urban counties-----	10	0	10	132	44	88	4,532	343
Rural and other urban areas-----	10	0	10	152	8	144	4,704	309
Northeastern Region-----	10	3	7	64	34	30	4,462	697
Largest SMSA's-----	3	3	0	3	3	0	5,013	5,013
Other large SMSA's----	2	0	2	5	5	0	4,589	1,836
Other SMSA's and highly urban counties-----	3	0	3	27	20	7	3,762	418
Rural and other urban areas-----	2	0	2	29	6	23	4,558	314
Midwestern Region-----	10	2	8	88	24	64	4,688	533
Largest SMSA's-----	2	2	0	2	2	0	5,279	5,279
Other large SMSA's----	3	0	3	18	18	0	4,604	767
Other SMSA's and highly urban counties-----	2	0	2	27	4	23	4,733	351
Rural and other urban areas-----	3	0	3	41	0	41	4,349	318
Southern Region---	10	0	10	116	29	87	4,297	364
Largest SMSA's-----	2	0	2	7	7	0	4,024	1,150
Other large SMSA's----	3	0	3	31	21	10	3,736	362
Other SMSA's and highly urban counties-----	3	0	3	48	0	48	4,891	306
Rural and other urban areas-----	2	0	2	30	1	29	4,519	301
Western Region----	10	2	8	96	35	61	4,476	466
Largest SMSA's-----	3	2	1	4	4	0	3,514	2,636
Other large SMSA's----	2	0	2	10	10	0	4,244	849
Other SMSA's and highly urban counties-----	2	0	2	30	20	10	4,945	330
Rural and other urban areas-----	3	0	3	52	1	51	5,280	305

<sup>1</sup>This total is larger than the 357 strata mentioned in the text. One reason for the difference is that several of the HIS self-representing strata were subdivided in designing the Cycle II sample. In addition, since the HIS sample was designed, two self-representing PSU's were split to form four PSU's, and one very small PSU which was omitted from the frame when the sample was drawn originally is designated "self-representing." Thus, there are actually 360 PSU's in the HIS design instead of 357.

The FSU's in the lowest density class were almost entirely non-self-representing. Although the average size of these FSU's was more than 300,000 in each region, each contained a number of PSU's. In fact, about 70 percent of all PSU's were classed as rural or small urban areas. These PSU's were quite small, typically containing only a few thousand people.

Two other modes of classification called "control classes" were added at the selection stage—rate of population change between 1950 and 1960 and geographic dispersion within regions, referred to as State groups.

The explicit use of the rate of population change is considered to be a major improvement in the design. It seems reasonable to view the rate of population change as a gross economic indicator and, consequently, a valuable health indicator. A depressed area can be generally characterized as having a below-average population gain and often a loss, whereas a new suburban area or new industrial area usually shows a large population increase.

Table D. Definition of rate-of-population-change classes by geographic region, 1950-60

Region	Rate of population change			
	$\alpha$	$\beta$	$\gamma$	$\delta$
	Percentage change			
Northeastern:				
SMSA PSU's--	<11	11-20	21 <sup>1</sup>	>21
Non-SMSA PSU's-----	<9	9-16	-	>16
Midwestern----	<6	6-18	19-25	>25
Southern-----	<5	5-21	22-42	>42
Western-----	<14	14-37	38-80	>80

<sup>1</sup>In the Northeastern Region, the two stands making up the New York SMSA constituted an entire rate-of-population-change class, giving it a single-value definition, a 21% increase.

Table E. State groups by geographic region

Region	State group
North-eastern----	1. Connecticut, Maine, Massachusetts, New Hampshire, Rhode Island, Vermont 2. New York 3. New Jersey, Pennsylvania
Midwestern--	1. Ohio 2. Indiana, Michigan, Wisconsin 3. Illinois 4. Minnesota 5. Iowa, Missouri
Southern----	1. Delaware, District of Columbia, Maryland, Virginia 2. Kentucky, Tennessee, West Virginia 3. Alabama, Arkansas, Louisiana, Mississippi 4. Georgia, North Carolina, South Carolina 5. Florida
Western-----	1. California 2. Oregon, Washington 3. Texas 4. Arizona, Colorado, Idaho, Montana, Nevada, New Mexico, Oklahoma, Utah, Wyoming, Alaska, Hawaii 5. Kansas, Nebraska, North Dakota, South Dakota

The rate-of-population-change classes also were defined on a sliding scale for each region, as indicated in table D, in such a way that each class contained approximately one-fourth of a region's population in 1960. Rate-of-population-change classes were defined slightly differently for SMSA's and non-SMSA's in the Northeastern Region. For other regions, no distinction was made between the two groups.

State groups within regions were instituted to maximize the spread of the sample among the States (table E). The basic criteria for forming State groups were to make the group membership as homogeneous as possible with respect to socioeconomic characteristics.



## SELECTION OF FIRST-STAGE UNITS

In addition to utilizing the fairly extensive stratification procedures described in the previous section, selection of units at the first stage of sampling also incorporated a modification of the Goodman-Kish controlled selection technique. This procedure permits some element of subjective determination in obtaining a "better balanced" or "more representative" sample, while retaining all the elements of true probability sampling. In particular, as used in this survey, it permitted proportional representation of the universe in several classes from each of five dimensions of classifications, even though only a grand total of 40 PSU's were selected.

The units sampled at the first stage of the HES sampling process were HIS strata. The term "first-stage unit" is employed to emphasize that, conceptually, the units being sampled were the aggregates of all PSU's in an HIS stratum. For example, in table F, the HIS sample PSU, Belknap-Merrimack, N.H., refers to seven PSU's constituting a single HIS stratum. This PSU with a population of 97,000 was the single PSU selected from among seven for the Health Interview Survey. The first of the FSU's from which a sample was selected in HES stratum Dii was the group of seven PSU's so referenced.

Prior to selecting the Cycle II FSU's, stratification was achieved for four broad population density classes within four geographic regions. As mentioned previously, this stratification resulted in a total of 40 HES superstrata—10 within each of the four geographic regions. Deeper stratification was precluded because of the requirement of selection of only one FSU from each superstratum. Had controlled selection not been used, and with no other restrictions except sampling with probability proportional to size, it would have been entirely possible, and indeed not improbable, that almost all the 10 sample PSU's in the Northeastern Region would be found to lie in the large metropolitan areas of New York, New Jersey, and Pennsylvania, with no representation at all from less populated areas such as Maine, Vermont, and New Hampshire.

An adaptation of the Goodman-Kish controlled selection technique was utilized which provided

for the identification of "control classes," constructed from variables other than the stratification variables, which were then used to reduce or eliminate such "batching" or extreme clustering of sample elements. Kish aptly refers to this as the introduction of "controls beyond stratification."<sup>22</sup> In the preceding example, the introduction of such a control would be used to increase the probability of inclusion of at least one FSU from States such as the three smaller ones named above, while maintaining the selection of one FSU from each stratum.

To the extent that the procedure is skillfully done, sampling variance is reduced. (Reduction is not certain and sampling variance may actually be increased.) Algebraic formulation of the impact of the procedure on sampling variances is not possible (or at least cannot be estimated from sample data from a single survey), but it is reflected in the half-sample replicate method generally used to estimate variances in this survey.<sup>23-25</sup> The control of probabilities by controlled selection is analogous to the formulation of balanced orthogonal patterns using Graeco-Latin squares familiar in experimental design, the major difference arising in increased complexity in calculating probability selection patterns due to unequal probabilities in the strata—control class cells. A good summary account of the fundamental concepts of controlled selection is given in reference 22, pages 488-495, and more detail may be found in the 1950 original article by L. Kish and R. Goodman.<sup>26</sup> The following discussion of the technique will be in the context of its application to the selection of the 10 FSU's for the Northeastern Region of the United States.

Classification of the 64 PSU's in the Northeastern Region into 10 superstrata on the basis of population density and size of FSU has been previously described. The superstrata are designated as Ai, Aii, Aiii, Bi, Bii, Ci, Cii, Ciii, Di, and Dii, where A indicates highest population density class, D the lowest, and i denotes the largest FSU-size class (iii the lowest).

Control classes were next defined using two additional variables—States group and rate of population change (from 1950-60). These classi-

fications for the Northeastern Region were as follows:

State group	Composition
(1)	Maine, New Hampshire, Vermont, Massachusetts, Rhode Island, Connecticut
(2)	New York
(3)	New Jersey, Pennsylvania

Rate of population change	Definition	
	SMSA PSU's	Other PSU's
$\alpha$	Under 11%	Under 9%
$\beta$	11-20%	9-16%
$\gamma$	21%	(empty cell)
$\delta$	22% and over	17% and over

Table F. Expected numbers of first-stage units and related data—HES superstratum Dii, Northeastern Region, Cycle II

State group	Rate of population change		HIS sample PSU (each representing one HIS stratum)	1960 Census of Population		Expected number	
	Class	Percent		HIS stratum	Control class		
1	$\alpha$	{ 7	Belknap-Merrimack, N.H. Kennebec-Lincoln, Maine	471,000	} 935,000	.19	
		5		464,000			
1	$\delta$	{ 37	Fairfield-Litchfield, Conn. Middlesex-New Haven, Conn. Hartford-Tolland, Conn. Bristol-Norfolk-Plymouth, Mass. Kent-Newport-Providence-Washington, R.I.	185,000	} 1,739,000	.36	
		30		418,000			
		28		308,000			
		21		455,000			
2	$\delta$	{ 17	Kent-Newport-Providence-Washington, R.I.	373,000	}		
2	$\alpha$	6	Chautauqua, N.Y.	338,000	338,000	.07	
		$\beta$	15	Chemung, Tioga-Tompkins, N.Y.	321,000	321,000	.07
			$\delta$	25	Orange-Putnam, N.Y.	265,000	265,000
3	$\alpha$	{ 7	Lycoming, Pa. Lebanon-Schuylkill, Pa.	256,000	} 520,000	.11	
		-5		264,000			
	$\beta$	13	Mercer, Pa.	127,000	127,000	.02	
3	$\delta$	{ 56	Monmouth-Ocean, N.J. Cumberland-Cape May, N.J.	443,000	} 598,000	.13	
		23		155,000			
TOTAL-----				4,843,000		1.00	

These variables define the controls beyond stratification which will relate to each stratum. Since one FSU is to be drawn from each superstratum and since the selection is to be with probability proportional to population size, the first step in the procedure is to determine expected numbers of FSU's in each control group by relating the populations of the control groups to a proportionate base of one FSU. For Cycle II, table F shows data for superstratum Dii.

Corresponding calculations for each superstratum result in expected numbers of FSU's for the full table of superstrata by control classes. These form the basic selection matrix for controlled selection analogous to the Graeco-Latin square. The full matrix for Cycle II data is shown in table G.

Values in the selection matrix show the expected numbers of sample FSU's which will be selected within any given cell. If the expected number is 1.0, exactly one FSU corresponding to that cell will be selected, and if the expected value is exactly zero, there will be no sample FSU's corresponding to that cell. If the expected number is 0.m, the probability is 0.m that one FSU corresponding to that cell will be selected, and 1-0.m that no sample FSU's corresponding to that cell will be selected.

The marginal row totals ensure that exactly one FSU will be selected from each superstratum, and the marginal column totals reflect the control beyond stratification of the control classes. For example, for control class  $\beta$  (3) the probability is .78 that two FSU's will be se-

Table G. Selection matrix for Northeastern Region, Cycle II, rate-of-population-change class and State group

HES superstratum	$\alpha$			$\beta$			$\gamma$			$\delta$			Total
	1	2	3	1	2	3	1	2	3	1	2	3	
Ai								1.00					1.00
Aii								1.00					1.00
Aiii						1.00							1.00
Bi					.31	.41						.28	1.00
Bii	.52		.48										1.00
Ci			.16	.22	.17					.14	.31		1.00
Cii	.08		.17	.17		.21			.22		.15		1.00
Ciii	.30		.27	.04		.14			.25				1.00
Di	.11	.19	.31	.11	.06						.09	.13	1.00
Dii	.19	.07	.11		.07	.02				.36	.05	.13	1.00
Total	1.20	.26	1.50	.54	.61	1.78		2.00		.97	.45	.69	10.00

lected from this class and .22 that only one will be selected. It is impossible that this control class will not be represented by any sample FSU.

Next, a set of selection patterns is developed which meet the requirements of the probabilistic restrictions of the selection matrix. The process is conveniently illustrated by a small, hypothetical case of two strata and two control classes.

Stratum	Control class		Total
	1	2	
I	.4	.6	1.0
II	.7	.3	1.0
Total	1.1	.9	2.0

The marginal limitations of the columns in this case imply that patterns may be formed by selecting 1 or 2 elements from class 1 and 0 or 1 element from class 2. Thus, 3 patterns are possible, namely:

Stratum	Control class	Pattern		
		1	2	3
I	1	1	1	0
	2	0	0	1
II	1	0	1	1
	2	1	0	0

Calculation of the probabilities of occurrence of these patterns involves solving the following

equations in which  $P_i$  is the probability associated with the  $i^{\text{th}}$  pattern,

$$\begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = \begin{bmatrix} .4 \\ .6 \\ .7 \\ .3 \\ 1.0 \end{bmatrix}$$

The last row of the coefficient matrix reflects the requirement that the sum of the probabilities of all patterns equals 1. For this simple and very restricted example there is a unique solution since the rank of the coefficient matrix equals three. However, in more complicated cases, the solution is usually not unique, and the judgmental decisions made in choosing patterns with nonzero probabilities influence the effectiveness of the procedure in achieving reduction of sampling variability. As the number of control classes and strata increase, the complexity of forming the selection patterns and calculating their associated probabilities increases rapidly. Kish has presented a method of forming selection patterns by successive subtraction of cell probabilities,<sup>22, 26</sup> and Schnack has developed a computer routine whereby sets of patterns may be generated and the resulting equations may be solved for the associated probabilities.<sup>27</sup> (There is no unique scheme which is favored by all, or even a majority, of samplers.)

For the Northeastern Region, a set of 17 patterns formed a complete set; that is, a set with associated probabilities totaling 1. The first six of these are indicated in table H. A single pattern is next chosen with probability proportional to the probability of occurrence of the pattern by selecting a random number between 0 and 1. For the data shown above the random number was .34 and pattern 3 was used in the survey.

A final selection is necessary for those cells in the pattern which contain more than one FSU. For example, table H shows that one FSU is to be selected in stratum Dii, control class 1.

Table H. Partial coefficient matrix of the first six of 17 selection patterns, HES Cycle II, Northeastern Region

HES superstratum	State group	Rate of population change	Pattern					
			1	2	3	4	5	6
Ai-----	2	$\gamma$	1	1	1	1	1	1
Aii-----	2	$\gamma$	1	1	1	1	1	1
Aiii-----	3	$\beta$	1	1	1	1	1	1
Bi-----	2 3 3	$\beta$	0	1	0	0	0	1
		$\beta$	1	0	1	0	1	0
		$\delta$	0	0	0	1	0	0
Bii-----	1 3	$\alpha$	1	0	1	0	1	0
		$\alpha$	0	1	0	1	0	1
Ci-----	1 1 2 3	$\beta$	0	0	0	1	0	0
		$\delta$	0	0	0	0	0	1
		$\beta$	1	0	0	0	0	0
		$\delta$	0	1	0	0	1	0
Cii-----	1 1 3 3	$\alpha$	0	0	0	0	0	0
		$\beta$	0	0	0	1	0	0
		$\delta$	1	0	0	0	0	0
		$\alpha$	0	0	0	0	1	1
Ciii-----	1 1 1 3	$\beta$	0	0	1	1	0	0
		$\delta$	0	0	0	0	0	0
		$\delta$	0	1	0	0	0	1
		$\alpha$	1	0	0	0	1	0
Di-----	1 2 2 3	$\alpha$	0	0	0	0	0	0
		$\beta$	0	0	0	0	0	1
		$\alpha$	0	0	0	1	0	0
		$\delta$	0	0	1	0	0	0
Dii-----	1 1 2 2 3 3	$\alpha$	0	1	0	0	0	0
		$\delta$	1	0	1	0	1	0
		$\alpha$	0	0	0	0	0	0
		$\beta$	0	0	0	0	0	0
		$\delta$	0	0	0	0	0	0
		$\alpha$	0	0	0	0	0	1
Probability of pattern-----			.17	.14	.07	.13	.06	.11
Cumulative probability-----			.17	.31	.38	.51	.57	.68

Table F shows the five FSU's constituting this cell. The final sampling operation selects one of the five with probability proportional to population size. If we denote  $x_{sj}$  as the size of the  $j$ 'th control group in the  $s$ 'th HES stratum and add a third subscript for the  $p$ 'th FSU within this cell, this final selection is with probability

$$\frac{x_{sjp}}{x_{sj}}$$

Further if we denote

$P_j$  = probability of selection of the  $j$ 'th cell

$P_\lambda$  = probability of selection of the  $\lambda$ 'th pattern

$X_s$  = population size of the  $s$ 'th super-stratum

$P_{\lambda j} = P_\lambda$  for all patterns which include  $j$   
 $= 0$  for all other patterns,

clearly  $P_j = \sum_\lambda P_{\lambda j}$ . However, since the original cell probabilities in table G are consistent

with the  $P_\lambda$ , it is true that  $\sum P_{\lambda j} = \frac{X_{sj}}{X_s}$

Thus the final probability of selection of each sample FSU was  $\frac{X_{sjp}}{X_s}$ . That is, within each

stratum, the probability of selection of each FSU was proportional to its population size, this sampling procedure having been maintained while providing the controls beyond stratification to reduce the probability of highly unrepresentative combinations and, hence, to achieve a reduction in sampling variability. The FSU, or HIS stratum, having been thus selected, the PSU previously selected to represent the HIS stratum, was then selected with probability 1 for purposes of the Health Examination Cycle II Survey. However, the actual probability of selecting the PSU from an FSU was proportional to the size of the FSU. Consequently, the probability of selecting a PSU

$$\text{was } \frac{X_{sj}}{X_s} \frac{X_{sjp}}{X_{sj}} = \frac{X_{sjp}}{X_s}$$

## WITHIN PSU DESIGN

### Problems of Development

A first-stage sample of 40 PSU's and the use of two mobile examining centers would permit the examination of about 180 children in each sample PSU, or a total of about 7,200 examinees

over a 2-year period. The within PSU design focused on the problem of selecting a probability sample of 8,000 children aged 6 to 11, or 200 in each sample PSU under the assumption that 90 percent of the children would be examined.

In developing the within PSU design, several problems had to be considered. The first was how to construct the universe, or sampling frame, to assure that every person in the target population has a chance of being selected in the sample. Secondly, there was some concern during the early stages of planning that parents would be reluctant to let their children travel long distances for an examination. One-way distances of 20 to 50 miles could be expected frequently, and occasionally more than 50 miles, if the sample segments were randomly selected throughout the PSU's. Thus, for large SMSA's and other PSU's covering large geographic areas, an intermediate stage of selection needed to be developed. Other problems to be considered in the within PSU design were how to select a sample of segments, or clusters of eligible children, and how to select a sample of children to be examined.

### Coverage of the Universe

The problem of selecting a probability sample of individuals is necessarily a complex one because there is no single best frame from which to select the sample and yet ensure complete coverage of the universe; in this case, the non-institutional population aged 6 through 11 residing outside Indian reservations. First, it will be useful to consider that the universe can be divided into four quadrants shown in the table below. The building blocks are the 1960 Census Enumeration Districts, which are small, well defined areas of about 200 housing units into which the entire Nation was divided for the 1960 Census of Population. Each enumeration district (ED) can be allocated to one and only one of the four quadrants according to a set of rules established by the Bureau of the Census. Enumeration districts whose 1960 Census Listing Books contain a high proportion of locatable or usable addresses are judged to be in either Quadrants A or C. Other ED's, mostly those with R.F.D. route addresses, are assigned to Quadrants B or D. The assignment to Quadrant A/B or C/D is based upon whether or not the ED is in a juris-

diction which maintains lists of building permits which can serve as a sampling frame. The approximate distribution of ED's, and consequently any sample of households, among the four quadrants is shown in parentheses in the following table:

	<i>Usable addresses</i>	<i>Not- usable addresses</i>	<i>Both types</i>
All areas-	(0.67)	(0.33)	(1.00)
Building permit areas--	A (0.57)	C (0.28)	(0.85)
Nonpermit area-----	B (0.10)	D (0.05)	(0.15)

The total universe of children eligible for the Health Examination Survey can be divided into the following four subuniverses:

- I. Eligible children residing in housing units listed in the 1960 census in ED's defined as having usable addresses.
- II. Eligible children residing in housing units listed in the 1960 census in ED's defined as not having usable addresses.
- III. Eligible children residing in housing units missed in the 1960 census.
- IV. Eligible children residing in housing units built since the 1960 census.

A PSU can, and usually does, contain ED's in each of the four quadrants. Furthermore, ED's generally contain children from three subuniverses, either I, III, and IV or II, III, and IV. Note that Subuniverses I and II are mutually exclusive. Subject to some possible errors in the application of the methods, coverage was made of the total universe by the following methods:

1. *Quadrant A.*—Subuniverse I was represented in the survey by a sample of clusters of addresses called list segments from 1960 Census Listing Books. Subuniverse III was given representation by a sample of "supplemental blocks." Supplemental blocks are chunks of land area, often a city block. For the Health Examination Survey, one supplemental block of about 24 housing units was selected for every three list segments selected. A map of each supple-

mental block and the 1960 Census Listing Book for the ED from which each was drawn were given to an interviewer for listing about 2 months prior to the initial interview date for the sample PSU. Any housing units in the supplemental blocks built prior to April 1960 and not listed in the 1960 Census Listing Book for their ED's were added to the sample under the assumption that they had been missed in the census. Subuniverse IV received coverage from a sample of building permits issued since April 1960.

2. *Quadrant B.*—The methods used to ensure coverage for ED's in Quadrant B differ from Quadrant A only in that both Subuniverses III and IV were given representation by the sample of supplemental blocks. This was accomplished by including in the sample all housing units not in the 1960 Census Listing Book, not only those built before April 1960.

3. *Quadrant C.*—Representation was given to Subuniverses II and III by a sample of small area segments selected from ED's defined as not having usable addresses and to Subuniverse IV by a sample of building permits issued after April 1960. Any overlap between the two frames was resolved by an inquiry into the date of construction of housing units in sample area segments and a deletion of any constructed after April 1960.

4. *Quadrant D.*—Finally, since no building permit data were available for ED's of Quadrant D, the area segments provided coverage for Subuniverse IV as well as for Subuniverses II and III. Since only about 5 percent of the sample is drawn from ED's in Quadrant D and there is probably little new construction in these predominantly rural areas, it is unlikely that there would be any sizable contribution to the mean square error arising from this quadrant.

### Selection of Localities Within PSU's

This intermediate stage of selection was considered important in the early stages of the cycle because it minimized the burden on the children and their parents by reducing the distance that some would have to travel to the examining center. If long distance travel should be a problem, then the selection of localities within sample PSU's should tend to maximize

the response rate and also reduce the cost of the survey.

The basic axis in the definition of a locality was in terms of Census minor civil divisions. Thus it was typically a city, part of a city, village, town, county, or the nonurbanized part of a county. The ultimate goal for a locality was that it should contain at least 250 children aged 6-11, or an elementary school enrollment of 250 children, or an area containing at least 2,000 people according to the 1960 census. The selection of an intermediate stage sample was not done routinely, but it was done on a PSU-by-PSU basis after a review of the problem had been made by NCHS-Census working committee.

Intermediate samples were selected for six PSU's only—Ashtabula-Geauga Counties, Ohio; Columbia-Dutchess Counties, New York; and the Denver, Philadelphia, Los Angeles, and Boston SMSA's. The procedure was discarded after the 10th stand because it was found that clustering sample segments in two or three areas, sometimes distant from feasible sites for the examination center, created an adverse situation. Random sampling of segments with probability proportional to size without an intermediate stage of sampling concentrated the sample around population centers where feasible examination centers could be located. Furthermore, there was little if any evidence that distance from a sample person's home to the examining site affected the participation rate or that mothers were reluctant to have their children travel so far. Also, any reduction in cost that accrued by sampling locations was more than compensated for by increased design efficiency resulting from the elimination of a stage of sampling.

For those six PSU's where locality sampling was used, after division into localities, a sample of three was drawn with a probability proportional to their 1960 population of children 6-9 years of age. In an SMSA one of the localities was the central city of the SMSA, and it was selected with certainty. From the remaining localities, which numbered from four to nine for the PSU's subsampled, two others were selected.

In four of the six PSU's subsampled, the Lahari sampling technique was used.<sup>28</sup> The method may be described briefly as follows:

1. Let the localities in a PSU be represented by  $L_1, L_2, \dots, L_1, \dots, L_m$ , which have measures of size  $A_1, A_2, \dots, A_1, \dots, A_m$ —the total number of children 6-9 years of age in each locality according to the 1960 census.

2. Let  $A_0$  be a number not smaller than the sum of  $m$  largest measures of size in the PSU.

3. Select, without replacement, a simple random sample of  $m$  localities.

4. Choose a random number  $R_1$  in the interval  $1 \leq R_1 \leq A_0$ .

5. If  $R_1 \leq \sum_{j=1}^m A_j$ , use the sample of size  $m$  selected in (3). If not, repeat the procedure until the condition is satisfied.

Because of the desire to control on the geographic spread of the sample in the Los Angeles and Philadelphia SMSA's, controlled selection of localities was used. The procedure will not be described since it is basically the same as that described above for the selection of first-stage units. However, it may be instructive to know how the populations were classified prior to sample selection.

The Philadelphia SMSA extends over the city of Philadelphia and three counties in New Jersey and four counties in Pennsylvania. The sampling plan was to select Philadelphia with certainty and two of the counties with a probability proportional to their 1960 population. To maximize the representativeness of the sample with respect to its urban-rural characteristics and geographic spread, the counties were grouped into a control selection matrix according to three degrees of urbanization classes (over 90%, 70-80%, and less than 70% urban) and State of location. Four controlled selection patterns of two counties each were formed. Then one of the four patterns was chosen by a random procedure.

The population of the Los Angeles SMSA was greater than the maximum size that had been set for a single stand (5 million) but was smaller than the minimum size of a double stand (8 million). This was also true of the Chicago SMSA, which had a 1960 population slightly below that of Los Angeles. To achieve a balance for the two areas it was decided to select 32 segments from the Los Angeles SMSA and 28



from Chicago, which, when combined, was the equivalent of three HES stands.

The Los Angeles SMSA was divided into the city of Los Angeles and four other strata of approximately 1 million people each. This stratification was accomplished by ranking the Census Minor Civil Divisions by their 1960 population size and dividing the total into quartiles. In addition to the city, one census division was drawn from each of the strata, controlling on geography and four population density-income classes.

### Selection of Segments

As stated in the section on coverage of the universe, there were four types of segments. List and area types of segments were the principal ones since they covered the vast majority of the target population. The other two, permit and supplemental block segments, were quite small since they included only eligible children residing in housing units built since the 1960 census and those residing in housing units missed in the 1960 census.

With only three exceptions, 20 segments were selected within each sample PSU from the frame of 1960 ED's contained in the sample PSU. Commonly, this sample of segments consisted of a combination of the two types depending on the character of the Census Listing Book of Addresses (usable or not) in the ED from which the segment was selected. In addition about three permit and supplemental block segments per PSU were selected, averaging about 1.5 eligible children per segment.

The area and list segments contained an expected nine children aged 5 to 9 in 1960 or about 11 children aged 6 to 11 at the time of the survey. Since the number of eligible children in a housing unit was a variable, there was a chance that 20 segments (plus the permit and supplemental block segments) would not yield the desired minimum sample of about 180 children. To overcome this potentiality, two reserve segments were selected, in addition to the 20, for the first eight stands. It became apparent at that time that 20 segments were sufficient, and therefore the selection of reserve segments was discontinued. The experience using this procedure was on the whole

satisfactory as indicated in table J, which shows the numbers of segments, interviewed housing units, and eligible children in the sample, by PSU and type of segment.

The sample of segments was selected in two steps. First, each ED was assigned a measure of size equal to a rounded whole number resulting from a division by 9 of the number of children aged 5 to 9 in the ED at the time of the 1960 census. Then a sample of 20 ED's was selected (except for the first eight stands when 22 were chosen) with probabilities proportional to the measures of size assigned to the ED's. Each sample ED was subsequently divided into as many roughly equal-sized segments, either area or list segments, as there were measures of size. The final step in the process was a random selection of one segment from each enumeration district. The selection procedure may be illustrated by a hypothetical example.

In the 1960 Census of Population, suppose a PSU was divided into 500 enumeration districts containing an average of about 200 housing units each. The addresses of the housing units were often not well-defined street numbers, so "area segments" were selected from this PSU.

For each ED the number of children aged 5 to 9 was determined as shown in the following table. Also shown are the appropriate "measures of size" resulting by dividing by 9 the number of children aged 5 to 9 in each ED, and the accumulation of measures of size over the entire PSU.

ED number	Number of children aged 5 to 9	HES measure of size	Accumulative measure of size
1	25	3	3
2	37	4	7
3	20	2	9
4	64	7	16
5	15	2	18
.	.	.	.
.	.	.	.
.	.	.	.
499	40	4	1,647
500	30	3	1,650
Total	15,000	1,650	1,650

Table J. Numbers of segments, interviewed housing units, and eligible children in the sample, by PSU and type of segment

PSU	Total			List and area segments			Permit and supplemental block segments		
	Segments	Interviewed housing units	Eligible children	Segments	Interviewed housing units	Eligible children	Segments	Interviewed housing units	Eligible children
Total-	954	21,393	8,589	820	20,928	8,382	134	465	207
1-----	28	630	200	22	615	195	6	15	5
2-----	25	475	246	22	469	241	3	6	5
3-----	26	638	248	22	628	239	4	10	9
4-----	23	602	218	22	602	218	1	0	0
5-----	25	600	230	22	592	227	3	8	3
6-----	25	459	206	22	448	204	3	11	2
7-----	31	505	240	22	473	224	9	32	16
8-----	26	451	240	22	446	236	4	5	4
9-----	22	410	248	20	402	240	2	8	8
10-----	20	727	147	16	708	143	5	19	4
11-----	24	777	201	20	740	191	4	37	10
12-----	24	694	138	16	679	130	8	15	8
13-----	24	546	246	20	520	234	4	26	12
14-----	23	459	196	20	439	188	3	20	8
15-----	22	539	193	20	534	193	2	5	0
16-----	22	882	220	20	882	220	2	0	0
17-----	23	689	195	20	676	193	3	13	2
18-----	23	395	241	20	387	239	3	8	2
19-----	24	727	226	20	708	221	4	19	5
20-----	24	423	252	20	410	242	4	13	10
21-----	21	379	218	20	373	217	1	6	1
22-----	21	495	234	20	493	233	1	2	1
23-----	37	690	301	32	673	288	5	17	13
24-----	23	451	160	20	442	156	3	9	4
25-----	20	434	221	20	434	221	0	0	0
26-----	25	408	188	20	402	183	2	6	5
27-----	22	338	186	20	330	184	2	8	2
28-----	22	267	179	20	263	176	2	4	3
29-----	25	528	239	20	507	231	5	21	8
30-----	23	421	149	20	400	139	3	21	10
31-----	24	450	216	20	437	207	4	13	9
32-----	24	506	250	20	498	248	4	8	2
33-----	25	650	260	20	626	247	5	24	13
34-----	20	422	239	20	422	239	0	0	0
35-----	23	680	231	20	675	226	3	5	5
36-----	24	492	218	20	478	209	4	14	9
37-----	22	596	222	20	589	220	4	7	2
38-----	26	616	228	20	595	226	6	21	2
39-----	22	545	163	20	540	159	2	5	4
40-----	21	397	156	20	393	155	1	4	1

To determine a sample of 20 ED's, divide 1,650 by 20 to get the length of the sampling interval (82.5). Select a random number between 0.1 and 82.5, say 13.0. This chooses ED number 4. The remaining 19 ED's are determined by adding this random number to the sampling interval and accumulating until the total exceeds 1,650 measures of size.

One segment is then selected from each sample ED. ED number 4 has a measure of size of 7; that is, it contains 7 segments with an average number of 11 children expected in each. Since these are area segments, it is necessary to identify their boundaries and approximate numbers of housing units contained within the segments. After the boundaries of the 7 segments have been determined, one of them is chosen, each having an equal probability of selection.

### Selection of Sample Children

The next step in the sample design was to identify the sample of children who were eligible to participate in the health examination. At each of the sample households a Census interviewer made a visit and asked certain questions. The questionnaire used is shown in appendix III.

The front of the questionnaire is concerned primarily with standard Census identification entries related to the housing unit. On the inside, the first group of questions that was asked identified the household composition. If there were no eligible children in the household, the interview was concluded with a few questions related to the possible presence of another household on the premises. In households in which there were eligible children, the remainder of the questionnaire was completed. A more detailed report of the administration of this questionnaire as well as the general plan, operation, and response results of the survey has been published.<sup>19</sup>

The 954 segments in the sample yielded a total of 25,106 households. Of these, 21,393 were interviewed, 2,291 were found to be vacant or to belong to persons having a usual residence elsewhere, and at 22 the composition of the household could not be established because of refusals or no one was home despite repeated calls.

In addition to the households identified above, 798 of the expected housing units in the original Census listing were found to have been demolished, outside segment boundaries, converted to business or storage, or merged with another unit.

The households interviewed yielded a total of 8,589 eligible children. The distribution of the number of segments, interviewed housing units, and eligible children for each PSU is shown in table J.

There was, however, a limit on the number of children that could be examined at a particular PSU. The time available for examinations at a particular PSU was necessarily set far in advance of any preliminary fieldwork. Therefore, the number of examinations that could be performed was dependent upon the number of examining days available. At most locations the number of days available, excluding Saturdays, was 18. The daily schedule of examinations called for six children in the morning and six in the afternoon so that 216 examining slots were available. However, because rescheduling was necessary for cancellations or no-shows, the maximum number of children who could be examined was approximately 200. At 26 locations, it was necessary to subsample the eligible children to yield around 190-200 sample children for examination.

Subsampling was accomplished through use of a master list which consisted of the names of eligible children determined in the household interviews. All eligible children in the PSU were listed in order by segment, serial (household order within segment), and column number (order in the household by age) and then numbered. After the desired subsampling rate had been determined, every  $n^{\text{th}}$  name on the list was deleted, starting with the  $y^{\text{th}}$  name,  $y$  being a number between 1 and  $n$  selected randomly. For example, if the total number of eligible children was 220, then a subsampling rate of one in 10 could be used which would reduce the number to 198. Selecting a random number between one and 10, say four, then the fourth eligible child on the master list would be deleted from the sample, as would every 10th following child, e.g., numbers 14, 24, 34, and 44.

## ESTIMATION PROCEDURE

An examination finding for an individual sample child is shown in data tabulations as a weighted frequency. This weight is a product of the reciprocal of the probability of selecting the child, an adjustment for nonresponse (not examined), and a poststratified ratio adjustment. The last was used to increase precision by bringing survey results into closer alignment with known U.S. population figures by color and sex within single years of ages 6 through 11.

The sample of slightly more than 7,400 children was arrived at by three stages of selection. The probability of an individual's being selected was the product of the probabilities of selection at three stages. In the first stage a single PSU was selected from each stratum. Within each sample PSU, a sample of segments expected to contain about 11 eligible children was selected. Then a subsample of the eligible was selected when the number of eligible children exceeded 200 in a sample PSU.

Since the strata are roughly equal in population size and a nearly equal number of sample children were examined in each of the sample PSU's, the sample design is essentially self-weighting with respect to the target population; that is, each child 6 to 11 years old has about the same probability of being drawn into the sample.

The adjustment for nonresponse is intended to minimize the impact of nonresponse on final estimates by imputing to nonrespondents the characteristics of "similar" respondents, that is, by relating nonrespondents to respondents by ancillary data known for both. Nonresponse due to refusals to be interviewed and "not-at-homes" amounted to only 22 households, so that the only nonresponse category requiring some adjustment was the "failure to be examined" nonresponses which amounted to 3.9 percent of the 7,417 sample children. "Similar" respondents were judged to be children in the same sample PSU having the same age (in years) and sex as the children not examined in the sample PSU. The weights of all respondents in a PSU having the same age and sex were adjusted upward to give representation to the nonrespondents in the PSU having that age and sex. Table K shows the total number of eligible children identified, the number of sample children, and the percent of sample children examined, by age and sex. The percent examined was quite similar for both boys and girls and for each age group. The response rate was also stable for each PSU ranging only from 90.6 to 100.0 percent as shown in table L.

The poststratified ratio adjustment used in the second cycle achieved most of the gains in precision which would have been attained if the sample had been drawn from a population stratified by age, color, and sex. The effect is to

Table K. Number of eligible children in the sample, number selected for examination, and percent examined, by age and sex

Age	Number of eligible children			Number of sample children			Percent of sample children examined		
	Total	Boys	Girls	Total	Boys	Girls	Total	Boys	Girls
Total-----	8,589	4,368	4,221	7,417	3,765	3,652	96.0	96.5	95.5
6 years-----	1,350	690	660	1,161	596	565	95.7	96.5	94.2
7 years-----	1,500	768	732	1,293	655	638	96.0	96.5	95.5
8 years-----	1,492	754	738	1,281	649	632	96.1	95.2	97.0
9 years-----	1,430	715	715	1,231	618	613	96.2	97.6	94.8
10 years-----	1,392	693	699	1,208	594	614	96.0	97.0	95.1
11 years-----	1,425	748	677	1,243	653	590	95.9	96.2	95.6

Table L. Number of sample children and number and percent examined, by stand number and location: Health Examination Survey, 1963-65

Stand location <sup>1</sup>	Stand number	Number of sample children	Examined	
			Number	Percent
All stands-----	...	7,417	7,119	96.0
Portland, Maine-----	1	200	198	99.0
Ashtabula, Ohio-----	2	185	175	94.6
Poughkeepsie, New York-----	3	193	190	98.4
Ottumwa, Iowa-----	4	196	195	99.5
Boston, Massachusetts-----	5	192	174	90.6
Denver, Colorado-----	6	192	189	98.4
Philadelphia, Pennsylvania-----	7	192	174	90.6
Lamar, Colorado-----	8	183	183	100.0
Charleston, South Carolina-----	9	186	171	91.9
Los Angeles, California-----	10 & 12	285	266	93.0
Sarasota, Florida-----	11	188	185	98.4
Atlanta, Georgia-----	13	191	187	97.9
San Francisco, California-----	14	189	187	98.9
Baltimore, Maryland-----	15	193	186	96.4
Mariposa, California-----	16	188	186	98.9
New York, New York-----	17 & 19	421	390	92.6
Moses Lake, Washington-----	18	193	189	97.9
Minneapolis, Minnesota-----	20	201	194	96.5
Grand Rapids, Michigan-----	21	191	186	97.4
Neillsville, Wisconsin-----	22	201	201	100.0
Chicago, Illinois-----	23	301	283	94.0
Des Moines, Iowa-----	24	160	159	99.4
Barbourville, Kentucky-----	25	196	185	94.4
Wichita, Kansas-----	26	188	178	94.7
Marked Tree, Arkansas-----	27	186	182	97.8
Brownsville, Texas-----	28	179	175	97.8
Houston, Texas-----	29	186	181	97.3
Birmingham, Alabama-----	30	149	144	96.6
Detroit, Michigan-----	31	168	162	96.4
Lapeer and Marysville, Michigan-----	32	179	175	97.8
Cleveland, Ohio-----	33	175	166	94.9
West Liberty and Beattyville, Kentucky-----	34	172	160	93.0
Allentown, Pennsylvania-----	35	173	159	91.9
Manchester and Bristol, Connecticut-----	36	174	167	96.0
Newark, New Jersey-----	37	177	167	94.4
Jersey City, New Jersey-----	38	175	163	93.1
Georgetown, Delaware-----	39	163	159	95.5
Columbia, South Carolina-----	40	156	148	94.9

<sup>1</sup>Cities in which trailers were located. Sample areas consisted of the PSU's which may have included several counties.

NOTE: Sample "take" for Los Angeles was deliberately somewhat low for "two stand locations" because that area should be only slightly over 1-1/2 stands on a population basis. Chicago, on the other hand, was oversampled in comparison with other "one stand locations," since it should be represented by slightly under 1-1/2 stands.

make the final sample estimates of population agree exactly with independent controls prepared by the Bureau of the Census for the U.S. non-institutional population at the midsurvey period (August 1, 1964) by color and sex for each single year of ages 6 through 11. The weights of every responding sample child in each of the 24 age, color, and sex classes were adjusted upward or downward so that the weighted total within the class equaled the independent population control. The poststratified adjustments required are shown in table M.

Table M. Poststratified adjustment factors (ratio of Census population control totals to Cycle II weighted estimates)

Age	White		All other	
	Boys	Girls	Boys	Girls
6 years-----	1.06	1.08	1.14	1.29
7 years-----	0.92	0.99	1.20	1.01
8 years-----	0.95	1.00	1.21	0.82
9 years-----	1.01	0.98	1.20	1.01
10 years-----	1.00	0.93	1.34	1.14
11 years-----	0.91	1.01	1.01	0.96

To aid in understanding the estimation procedure, the estimator is presented as follows:

Consider an  $X$ -characteristic of the  $l^{\text{th}}$  sample person in the  $k^{\text{th}}$  segment,  $j^{\text{th}}$  age-sex class,  $i^{\text{th}}$  PSU,  $h^{\text{th}}$  stratum, and the  $g^{\text{th}}$  age-sex-color class in the United States, denoted by  $X_{g,h,i,j,k,l}$ . An estimator,  $X'$ , of a total aggregate,  $X$ , in the U.S. population is derived from Cycle II data using the following equation:

$$X' = \sum_{g=1}^{24} R'_g \sum_{h=1}^{40} W_{1,h} \sum_{i=1}^{40} W_{2,hi} \sum_{j=1}^{12} \frac{n'_{.hij}}{n_{.hij}}$$

$$\sum_{k=1}^s W_{3,hi,k} \sum_{l=1}^{\dot{n}_{.hij,k}} X_{g,hij,k,l}$$

where  $R'_g = Y_g / Y'_g$  = ratio of total U.S. non-institutional population in the  $g^{\text{th}}$  age-sex-color group according to the 1964 census figures to the estimated total U.S. noninstitutional

population in the  $g^{\text{th}}$  age-sex-color group using a simple inflation-type estimator, adjusted for nonresponse.

$W_{1,h}$  = first-stage design weight for the  $h^{\text{th}}$  stratum (i.e., superstratum) =  $\frac{1}{P_{1h}}$  the reciprocal of the probability of selecting a PSU from the  $h^{\text{th}}$  stratum.

$W_{2,hi}$  = second-stage design weight for the  $i^{\text{th}}$  PSU in the  $h^{\text{th}}$  stratum =  $\frac{1}{P_{2hi}}$  the reciprocal of the probability of selecting a segment from the  $i^{\text{th}}$  PSU in the  $h^{\text{th}}$  stratum.

$W_{3,hi,k}$  = third-stage design weight = the reciprocal of the probability of selecting a person in the subsample from eligible persons in the  $k^{\text{th}}$  segment,  $i^{\text{th}}$  PSU,  $h^{\text{th}}$  stratum.

$n'_{.hij}$  = total eligible persons, after subsampling, in the  $j^{\text{th}}$  age-sex class in the  $i^{\text{th}}$  PSU and  $h^{\text{th}}$  stratum.

$\dot{n}_{.hij}$  = total examined eligible persons in the  $j^{\text{th}}$  age-sex class in the  $i^{\text{th}}$  PSU and  $h^{\text{th}}$  stratum.

In addition to the adjustment factors indicated in the equation, another adjustment of  $\frac{22}{20}$  was applied to data collected in the first eight stands completed since 22 "regular" segments per PSU were originally selected and only 20 were used. The distribution of final estimation weights is shown in table N.

Table N. Distribution of final estimation weights for examined children

Weight class	Distribution of examined children	
	Number	Percent
1,000-1,999-----	205	2.9
2,000-2,999-----	2,762	38.8
3,000-3,999-----	3,040	42.7
4,000-4,999-----	837	11.8
5,000-5,999-----	82	1.1
6,000-6,999-----	64	0.9
7,000-8,999-----	58	0.8
9,000-10,999-----	44	0.6
11,000-14,279-----	27	0.4

## VARIANCE ESTIMATION

### Background

Standard errors of estimates of parameters for the sample were estimated by means of the (balanced) half-sample replication technique, first adapted to use for large-scale surveys by Simmons and Losee, described in references 23-25, 29-31. The reasons for the adoption of this method were both operational and theoretical. The following major characteristics of the survey suggested requirements that were largely or wholly met by the half-sample replication technique.

1. Since the obtaining of data for any single sample child is costly, the sample size is necessarily limited. The obvious statistical objective of maximum exploitation of the data is particularly meaningful in the context of the Health Examination Survey since an increase in sample size has an immediate and consequential impact on costs. The Health Examination Survey cannot afford, for example, overuse of commonly employed "upper limit" approximations to sampling errors as might be done with a large sample group.

2. Because the sampling errors of most statistics are large enough to be meaningful in analysis and many are large enough to be critical to the analytical conclusions, a high degree of computational support for the researchers analyzing the material is indicated. Standard errors must be made available quickly so that a particular investigation, which frequently advances in stepwise fashion with the next analytical step depending on the results of the last, may proceed with reasonable speed.

3. The complete algebraic formula for estimation of sampling errors for the survey design is unknown. This is because of the nature and complexity of the design as described in the preceding sections. While the algebraic relationships are identifiable or capable of being developed for particular subprocedures—such as the use of cluster and multistage sampling within strata to reduce costs, the poststratification techniques used to reduce sampling variance, or the nonresponse adjustments to reduce bias—a single, composite, estimating equation for the standard error of survey statistics

cannot be developed. The use of the Goodman-Kish controlled selection technique as part of the selection process in itself precludes this, since, while it is known that such controlled selection should reduce the sampling variance, theory does not exist to permit algebraic quantification of the extent of the reduction using only sample information. Even if controlled selection were eliminated as a definitive factor, the extreme complexity of the combination of the various other elements of the design would probably preclude, as a matter of practicality, direct algebraic estimation.

4. In a large, multidimensional investigation, such as the National Health Examination Survey of children, interest frequently centers on studies of characteristics of various population subgroups. The numbers of persons in these subgroups, or domains of study, are in themselves random variables. Algebraic techniques for computation of standard errors of statistics relating to them have been developed by Cochran<sup>32</sup> and others for certain restricted designs, all considerably less involved than the survey design used for Cycle II. Their use, however, introduces some bias, considerable complexity, and formidable computational effort.

### Summary of Applicable Theory

The population is classified into  $L$  strata, from each of which two sample PSU's are drawn, with equal probability within the stratum, but not necessarily across strata. The desideratum of selection of exactly two sample PSU's reflects an essential element of the theory and may be met by post facto "collapsing" of two strata from each of which only a single PSU has been drawn or by creating an artificial PSU by random methods from the operational PSU selected from such strata.

Of analytical interest is parameter  $P$  for which an estimate,  $p$ , is to be obtained from the sample. The estimator,  $p$ , is a linear combination of the sample observations in fully rigorous developments, although, as will be seen later, this requirement may be compromised in applications with little practical effect.

A half-sample replicate is defined as the collection of  $L/2$  PSU's obtained by selecting one of the paired sample PSU's from each stratum. (These may be referred to simply as replicates or half samples for brevity.) Designating  $i=1, 2$  as the subscript to identify the sample PSU's within each stratum,  $h=1, 2, 3, \dots, L$  to identify the strata, and  $a=1, 2, 3, \dots, A$ , where  $A \geq L$  as half-sample replicate identification, the pattern may be summarized as in table O where a "+" indicates that a PSU falls into the particular half sample, and a "-" indicates that it does not.

Analogues of the linear estimator  $p$  corresponding to each half sample are then computed. That is, for the  $a$ 'th half sample,  $P_a$  is calculated by summing across strata as:

$$P_a = \sum_{h=1}^L w_h P_{hi}$$

where  $w_h$  is the proportion of persons in the  $h$ 'th stratum ( $\sum w_h = 1$ ),  $i$  is either 1 or 2 depending on which PSU of the stratum is in

Table O. Half-sample replication formation

Half-sample replication	Stratum									
	1		2		3		...	L		
	PSU		PSU		PSU			PSU		
	1	2	1	2	1	2	1	2		
1	+	-	-	+	-	+	...	+	-	
2	-	+	-	+	+	-	...	-	+	
3	-	+	+	-	-	+	...	-	+	
.	.	.	.	.	.	.	...	.	.	
.	.	.	.	.	.	.	...	.	.	
.	.	.	.	.	.	.	...	.	.	
A	+	-	+	-	+	-	...	+	-	

half-sample  $a$  and  $P_{hi}$  is, in this example, a mean.

The estimator  $p$  calculated using all the information in the sample, is:

$$p = \sum_h \frac{w_h (P_{h1} + P_{h2})}{2}$$

The variance of the estimator  $p$  is calculated as:

$$S_p^2 = \frac{1}{A} \sum_{a=1}^A (P_a - p)^2$$

A set of side conditions relating to the selection of PSU's for development of the half samples has been developed by McCarthy,<sup>24, 25</sup> based on work by Plackett-Burman<sup>33</sup> and Gurney.<sup>34</sup> The significance of this procedure is that greatly increased stability in the estimate  $S_p^2$  is obtained by eliminating a between-strata contribution of variance otherwise present in calculating  $S_p^2$  across half samples. The  $S_p^2$  calculated for a set of half samples formed according to the McCarthy criteria is numerically equal to the value which would be obtained if all  $2^L$  possible half samples had been formed. A set of half samples selected according to the McCarthy criteria is called a *balanced* set and the procedure is referred to as balanced half-sample pseudoreplication. The Cycle II variance estimations are calculated by using balanced half-sample replication methods, and reference to the technique throughout this report implies a balanced pattern.

Estimates of standard errors developed according to this technique have several highly desirable attributes, both in calculation and in concept. The more important are summarized by McCarthy<sup>24</sup> as:

"Replicated sampling permits one to bypass the extremely complicated variance estimation formulas and the attendant heavy programming burdens. Variance estimates



based upon the replicated estimates will mirror the effects of all aspects of sampling and estimation that are permitted to vary randomly from replicate to replicate. This of course includes the troublesome domain-of-study problem."

The theory is completely rigorous only in the case in which the statistics for which standard errors are being estimated are linear functions of the sample observations. Several empirical investigations indicate that use of certain ratio estimators and correlation statistics results in a bias that is unimportant, if detectable at all, in an analytical context.<sup>23-25, 29</sup> Such bias is not considered to be of practical importance in application of the replication method to Cycle II data, as described below.

### Application to Cycle II Data

The starting point for Cycle II replication procedures is the set of 40 PSU's, one from each of the 40 HES superstrata as previously defined. Associated with each of these PSU's is a sampling fraction which is numerically equal to the probability of selection of the PSU although, as described in a preceding section, the actual mechanics of selection of the PSU involve application of the Goodman-Kish side conditions which are more complicated (and contribute more to reduction in sampling variation) than simple selection of the PSU with probability proportional to size. An example will clarify the way in which the weights associated with the sample PSU's were computed.

In the Northeastern Region, superstratum Ciii is composed of 11 HIS strata or FSU's with a combined 1960 census population of 3,759,516 (table P). This HES stratum consisted of SMSA's of under 1,000,000 population in 1960 (C designation), which contained the smaller SMSA's (iii designation) in this category.

HES superstratum Ciii includes HIS stratum No. 211 which is in turn composed of two HIS PSU's: Portland, Maine, SMSA (1960 census population 120,655) and Atlantic City, N.J., SMSA (1960 census population 160,880).

Under the Goodman-Kish selection technique, HIS stratum No. 211 is selected from the 11 HIS strata which constitute HES superstratum

Ciii. The Portland, Maine, HIS sample PSU which has already been drawn from HIS stratum No. 211 for HIS purposes is then selected for HES purposes with probability 1 and is designated as the HES "stand." The numerical value of the probability of selection of the Portland, Maine, stand in this case is:

$$\frac{120,655}{120,655 + 160,880} + \frac{120,655 + 160,880}{3,759,516}$$

although, as explained in a previous section, the actual (Goodman-Kish) selection procedure resulting in this probability is operationally different from simple probability proportional to size selection which might be (incorrectly) inferred from the above two fractions. The actual selection procedure is also conceptually different since the Goodman-Kish side conditions result in a smaller sampling variance.

The stands, or examination locations, corresponding to the PSU's thus selected are identified in table P together with the HES superstrata with which they are associated.

As stated previously, the balanced half-sample replication theory is based upon selection of one sampling unit from a stratum containing exactly two such units. It was therefore necessary at this point to create HES artificial or "pseudo" strata from pairs of HES strata in order to make use of the half-sample replication model. Two procedures were used, depending on whether or not the defined HES strata were self-representing.

For both self-representing (certainty) and non-self-representing (noncertainty) HES strata, strata were paired on the basis of (1) some subjective determination of the homogeneity of the population in which the primary considerations were population density, region, rate of growth, and industry and (2) concern that strata of approximately equal size would be paired. The latter has no theoretical or practical effect on variance computations in Cycle II since the factors necessary to adjust for unequal size of members of the pair were introduced into the weighting procedures specific for each replication (reference 22, page 285). The former is of concern, in that members of the pair may have markedly different characteristics with respect

Table P. Definition of HES pseudostrata for replication purposes

HES superstratum	1960 Census of Popula- tion of HES super- stratum	Region	HES pseudo- stratum number	Stand	
				Num- ber	Location
Non-self-representing HES strata					
Bii-----	4,994,736	NE	01	5	Boston, Mass.
Bi-----	4,183,250	NE	01	37	Newark, N.J.
Ci-----	3,759,760	NE	02	38	Jersey City, N.J.
Cii-----	3,768,466	NE	02	35	Allentown, Pa.
Di-----	4,271,826	NE	03	3	Columbia-Dutchess, N.Y. (Poughkeepsie, N.Y.)
Dii-----	4,843,253	NE	03	36	Hartford-Tolland, Conn. (Manchester and Bristol, Conn.)
Bii-----	3,776,544	S	04	40	Columbia, S.C.
Biii-----	3,961,447	S	04	9	Charleston, S.C.
Cii-----	4,961,779	S	05	27	Crittenden-Poinsett (Marked Tree, Ark.)
Di-----	4,622,338	S	05	39	Sussex, (Georgetown, Del.)
Ciii-----	4,973,857	S	06	25	Bell-Knox-Whitley, Ky. (Barbourville)
Dii-----	4,415,267	S	06	34	Breathitt-Lee, Ky. (West Liberty and Beattyville)
Bi-----	3,856,698	MW	07	33	Cleveland, Ohio
Bii-----	5,155,715	MW	07	20	Minneapolis-St. Paul, Minn.
Di-----	4,507,428	MW	08	32	Lapeer-St. Clair, Mich. (Lapeer and Marysville)
Dii-----	4,156,090	MW	08	2	Ashtabula-Geauga, Ohio
Aiii-----	3,890,572	W	09	14	San Francisco, Calif.
Bii-----	4,899,898	W	09	6	Denver, Colo.
Dii-----	5,519,588	W	10	8	Prowers, Colo. (Lamar)
Diii-----	5,115,227	W	10	16	Mariposa, Calif.
Aii-----	4,318,307	S	11	13	Atlanta, Ga.
Bi-----	3,587,125	W	11	29	Houston, Tex.
Ci-----	4,895,507	MW	12	24	Des Moines, Iowa
Ci-----	5,047,027	W	12	26	Wichita, Kans.
Bi-----	3,472,118	S	13	30	Birmingham, Ala.
Biii-----	4,799,314	MW	13	21	Grand Rapids, Mich.
Diii-----	4,384,792	MW	14	22	Clark, Wis. (Neillsville)
Di-----	5,207,020	W	14	18	Grant, Wash. (Moses Lake)
Ciii-----	3,759,516	N	15	1	Portland, Maine
Cii-----	4,570,419	MW	15	4	Mahaska-Wapello, Iowa (Ottumwa)
Ci-----	4,739,463	S	16	11	De Soto-Sarasota, Fla. (Sarasota)
Cii-----	4,841,990	W	16	28	Brownsville, Tex. (Brownsville)

Table P. Definition of HES pseudostrata for replication purposes—Con.

HES superstratum	1960 Census of Popula- tion of HES super- stratum	Region	HES pseudo- stratum number	Stand	
				Num- ber	Location
Non-self-representing HES strata					
Aiii Ai	4,342,897 3,728,920	NE S	01A,01B 01A,01B	7 & 15	Philadelphia, Pa., and Baltimore, Md.
Ai Aii	6,794,461 3,762,360	MW MW	02A,02B 02A,02B	23 & 31	Chicago, Ill., and Detroit, Mich.
Ai,Aii Ai,Aii	6,742,696	W	03A,03B 03A,03B	10 12	Los Angeles, Calif.
Ai,Aii Ai,Aii	10,694,633	NE	04A,04B 04A,04B	17 & 19	New York, N.Y.

to a particular variable under study. To the extent that this is true then the expected value of the estimated standard error may be positively biased to some extent. That is, as the subjective pooling of "collapsing" of strata becomes a compromising procedure, a more conservative estimate (i.e., overstatement) of the sampling variance is obtained (reference 22, page 283). Evaluation of this effect for Cycle II data suggests that any resulting overstatement of sampling variance is of trivial consequence in an analytical context.

The specific pairing or "collapsing" procedures used for Cycle II are indicated in table P.

For self-representing strata, an additional procedure was followed to ensure homogeneity of populations. This is best described in terms of an example using the first two self-representing superstrata identified in table P. After the pairing of HES superstrata Aiii (NE) and Ai(S), sample segments in the Philadelphia and Baltimore PSU's were selected in random serpentine fashion so that HES Pseudo-PSU 01A, the population corresponding to half of the segments, includes a randomly defined part of both the Philadelphia SMSA and Baltimore SMSA

populations. This is, of course, also true for HES Pseudo-PSU 01B. These two Pseudo-PSU's constitute HES Pseudostratum 01.

As indicated in table P, Los Angeles and New York are special cases in which a single HES pseudostratum was defined from a single HES superstratum, the usual procedure, of course, being the definition of a single HES pseudostratum from the two HES superstrata. They were, however, subjected to the randomization process described in the preceding paragraph, even though the artificially defined "stands" for these areas had already been defined on the basis of randomly selected segments with no geographical clustering.

For non-self-representing strata, the pseudostrata were defined on the basis of size and homogeneity of population as shown in table P.

Having defined the 20 (artificial) pseudostrata, each consisting of two PSU's, the balanced half-sample replication pattern following the Plackett-Burman techniques may be applied. This was done, and 20 half-sample replications were formed according to the constraints developed by Plackett-Burman. Each (half-sample) replication consisted of 20 sample PSU's, one being selected from each pseudostratum.

One additional ramification was undertaken before variance computations were made. This was the development and application of factors to adjust each individual replication to the (Census) independent control populations for 24 age-color-sex classes. For example, the combined sampling and nonresponse weights for 8-year-old white male children in replication four were adjusted so that the national estimate of all such children, using only the sample information contained in replication four, results in a figure of 1,739,000—the independent Census estimate of this population as of August 1, 1964. In summary, each replication (which contains about half of the sample cases) results in an estimate which is numerically equal to the estimate obtained from the whole sample due to the application of these adjustment factors. While this reduces a small amount of bias of the estimated sampling variance, the process involves considerable work and insufficient evidence is available on which to base a decision as to whether or not it is worth the cost.<sup>23,29</sup> Pending further methodological investigations, a prudent approach was adopted for Cycle II data and the factors were applied as described.

The only remaining step is the application of the theory stated earlier to produce the variance estimates. To avoid restatement of the theory, application will be noted in the form of an example, paralleling the theory presented earlier.

Data from Cycle II show that the mean number of upper arch permanent teeth among 8-year-old boys in families for which the annual family income is reported as between \$5,000 and \$6,999 is 5.17, i.e.,  $p=5.17$  using the previous notation. For each of the 20 half-sample replicates, the analogue  $p_a$  is computed (table Q). The sampling variance of  $p$  is then estimated as

$$S_p^2 = \frac{1}{20} \sum_{a=1}^{20} (P_a - p)^2$$

$$= .008545$$

For analytical convenience several functions of the estimated sampling variance are then calculated and routinely displayed. The values of

these for this estimate of the mean number of upper arch permanent teeth are as follows:

Mean upper arch permanent teeth-----	5.17
Standard error of mean-----	.09
Estimated population (denominator)-----	437,000
Standard error of denominator	34,000
Estimated upper arch permanent teeth (numerator)-----	2,258,000
Standard error of numerator--	178,000
Rel-variance of mean-----	.00032
Rel-variance of denominator--	.00666
Rel-variance of numerator----	.00625
Sample frequency-----	140

A standard computer program is available whereby means, standard errors of means, sample sizes, and the associated indexes of sampling variability are obtained for a cross-classification of about 300 cells with simple and routine specifications. Row percentages and rates with associated statistics are also options. Replicate variance calculations are also programmed for correlation and regression statistics, although at this writing, data processing restrictions limit use of this latter program to methodological investigations rather than for routine analytical purposes.

Table Q. Half-sample replicate estimates of mean number of upper arch permanent teeth for 8-year-old boys with family income of \$5,000-\$6,999

Replicate number	$P_a$	Replicate number	$P_a$
1	5.1029	11	5.1899
2	5.0685	12	5.0066
3	5.1964	13	5.2291
4	5.2701	14	5.2074
5	5.1602	15	5.0424
6	5.2353	16	5.0260
7	5.1779	17	5.2465
8	5.2547	18	5.3713
9	5.1619	19	5.1005
10	5.1116	20	5.0737

## REFERENCES

<sup>1</sup>National Center for Health Statistics: Origin, program, and operation of the U.S. National Health Survey. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 1-No. 1. Public Health Service. Washington. U.S. Government Printing Office. Apr. 1965. (reprint)

<sup>2</sup>National Center for Health Statistics: Plan and initial program of the Health Examination Survey. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 1-No. 4. Public Health Service. Washington. U.S. Government Printing Office. July 1965.

<sup>3</sup>U.S. National Health Survey: A study of special purpose medical-history techniques. *Health Statistics*. PHS Pub. No. 584-D1. Public Health Service. Washington. U.S. Government Printing Office. Jan. 1960.

<sup>4</sup>U.S. National Health Survey: Attitudes toward cooperation in a health examination survey. *Health Statistics*. PHS Pub. No. 584-D6. Public Health Service. Washington. U.S. Government Printing Office. July 1961.

<sup>5</sup>U.S. National Health Survey: Evaluation of a single-visit cardiovascular examination. *Health Statistics*. PHS Pub. No. 584-D7. Public Health Service. Washington. U.S. Government Printing Office. Dec. 1961.

<sup>6</sup>National Center for Health Statistics: Comparison of two vision-testing devices. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 1. Public Health Service. Washington. U.S. Government Printing Office. June 1963.

<sup>7</sup>National Center for Health Statistics: The one-hour oral glucose tolerance test. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 3. Public Health Service. Washington. U.S. Government Printing Office. July 1963.

<sup>8</sup>National Center for Health Statistics: Cooperation in health examination surveys. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 9. Public Health Service. Washington. U.S. Government Printing Office. July 1965.

<sup>9</sup>National Center for Health Statistics: Replication, an approach to the analysis of data from complex surveys. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 14. Public Health Service. Washington. U.S. Government Printing Office. Apr. 1966.

<sup>10</sup>National Center for Health Statistics: Three views of hypertension and heart disease. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 22. Public Health Service. Washington. U.S. Government Printing Office. Mar. 1967.

<sup>11</sup>National Center for Health Statistics: Factors related to response in a health examination survey. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 36. Public Health Service. Washington. U.S. Government Printing Office. Aug. 1969.

<sup>12</sup>National Center for Health Statistics: *Vital and Health Statistics*. PHS Pub. No. 1000-Series 11-Nos. 1-37. Public Health Service. Washington. U.S. Government Printing Office.

<sup>13</sup>National Center for Health Statistics: Evaluation of psychological measures used in the health examination survey of children ages 6-11. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 15. Public Health Service. Washington. U.S. Government Printing Office. Mar. 1966.

<sup>14</sup>National Center for Health Statistics: Calibration of two bicycle ergometers used by the Health Examination Survey. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 21. Public Health Service. Washington. U.S. Government Printing Office. Feb. 1967.

<sup>15</sup>National Center for Health Statistics: A study of the achievement test in the Health Examination Surveys of persons aged 6-17 years. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 24. Public Health Service. Washington. U.S. Government Printing Office. June 1967.

<sup>16</sup>National Center for Health Statistics: Orthodontic treatment priority index. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 25. Public Health Service. Washington. U.S. Government Printing Office. Dec. 1967.

<sup>17</sup>National Center for Health Statistics: Development of the brief test of literacy. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 27. Public Health Service. Washington. U.S. Government Printing Office. Mar. 1968.

<sup>18</sup>National Center for Health Statistics: *Vital and Health Statistics*. PHS Pub. No. 1000-Series 11-Nos. 101-103. Public Health Service. Washington. U.S. Government Printing Office.

<sup>19</sup>National Center for Health Statistics: Plan, operation, and response results of a program of children's examination. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 1-No. 5. Public Health Service. Washington. U.S. Government Printing Office. Oct. 1967.

<sup>20</sup>National Center for Health Statistics: The Statistical Design of the Health Household-Interview Survey. *Health Statistics*. PHS Pub. No. 584-A2. Public Health Service. Washington. U.S. Government Printing Office. July 1958.

<sup>21</sup>Bureau of the Census: *The Current Population Survey*. A Report on Methodology. Technical Paper No. 7. U.S. Department of Commerce. Washington. U.S. Government Printing Office. 1963.

<sup>22</sup>Kish, L.: *Survey Sampling*. New York. John Wiley and Sons, Inc., 1965.

<sup>23</sup>Simmons, W. R., and Baird, J.: Use of pseudoreplication in the NCHS health examination survey. *Proceedings of the Social Statistics Section of ASA*. American Statistical Association. 1968.

<sup>24</sup>National Center for Health Statistics: Replication, an approach to the analysis of data from complex surveys. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 14. Public Health Service. Washington. U.S. Government Printing Office. Apr. 1966.

<sup>25</sup>National Center for Health Statistics: Pseudoreplication, further evaluation and application of the balanced half-sample technique. *Vital and Health Statistics*. PHS Pub. No. 1000-Series 2-No. 31. Public Health Service. Washington. U.S. Government Printing Office. Jan. 1969.

<sup>26</sup>Goodman, R., and Kish, L.: Controlled selection—a technique in probability sampling. *J. Am. Statist. A.* 45:350-373. 1950.

<sup>27</sup>Schnack, G.A.: Personal Communication.

<sup>28</sup>Lahiri, D.B.: A method for sample selection providing unbiased ratio estimates. *International Statistical Institute Bulletins*. 33,2, 133-140. 1951.

<sup>29</sup>Kish, L. and Frankel: *Proceedings of the Social Statistics Section of the ASA*. American Statistical Association. 1968.

<sup>30</sup>Losee, G.J.: Experiments in estimation and calculation of variability for statistics from the Health Examination Survey. *Biometrics*. Dec. 1964.

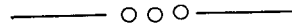
<sup>31</sup>Simmons, W.R., and Bean, J.A.: Impact of design and estimation components on inference. *New Developments in Survey Sampling*. Wiley-Interscience. 1969.

<sup>32</sup>Cochran, W. G.: *Sampling Techniques*. ed. 2. New York. John Wiley and Sons, Inc. 1963.

<sup>33</sup>Plackett, R. L., and Burman, J.P.: The design of optimum multifactorial experiments. *Biometrika*. (33):305-325. 1946.

<sup>34</sup>Gurney, M.: The Variance of the Replication Method for Estimating Variances for the CPS Sample Design. Dittod Memorandum. U.S. Bureau of the Census. 1962.

<sup>35</sup>Jessen, R.L.: Some methods of probability non-replacement sampling. *J. Am. Statist. A.* 64:325. 1969.



## APPENDIX I

### GLOSSARY OF TERMS

*Primary sampling unit (PSU).*—A geographic entity composed of one or more contiguous counties, or a standard metropolitan statistical area (SMSA). The 3,103 counties and independent cities in the United States were grouped to form 1,891 PSU's. Details of how PSU's were formed are presented in the text of this report.

*Self-representing PSU's.*—Those PSU's which cover an entire stratum. The 1,891 PSU's were grouped into 357 HIS strata. Of these, 112 are composed of a single PSU and 245 contain more than one. Since one PSU was selected in the sample from each stratum, those strata containing only one PSU are *self-representing* and those containing more than one PSU are *non-self-representing*.

*First-stage units (FSU's).*—With a few exceptions, an FSU is synonymous with an HIS stratum, consisting of the aggregate of PSU's, sample and non-sample, in the stratum.

*HES superstratum.*—Consists of one or more FSU's. For the Cycle II sample, 364 FSU's were grouped into 40 superstrata. Eight superstrata were self-representing, and 32 were non-self-representing.

*Pseudostratum.*—An artificial stratum formed by combining two superstrata, or by combining "random halves" of superstrata. Examples of the latter are two pseudostrata, each comprised about half of the population of the Philadelphia SMSA plus half of the population of the Baltimore SMSA. The pseudostrata were conceptual entities used in the estimation of variances by the half-sample replication method. Twenty pseudostrata were defined, 16 from the combining of two superstrata.

*Standard metropolitan statistical area (SMSA).*—A county or group of contiguous counties (except in New England) which contains at least one central city of 50,000 people or more, or "twin cities" with a combined population of at least 50,000 population. In addition, other contiguous counties are included in an SMSA if, according to certain criteria, they are socially and economically integrated with the central city. A detailed explanation of a listing of the component areas of each SMSA is given in Bureau of the Budget, *Standard Metropolitan Statistical Areas*, 1967 Edition.



## APPENDIX II

### PROCEDURE FOR FORMING AND STRATIFYING PSU'S IN THE CURRENT POPULATION SURVEY AND THE HEALTH INTERVIEW SURVEY DESIGNS

#### Formation of PSU's

Several rules were followed in defining and forming PSU's. They were:

1. Each PSU should comprise one or more contiguous counties. PSU's involving metropolitan counties were defined as consisting of whole SMSA's, except in New England where towns and cities rather than counties were used in defining SMSA's. (For definition of SMSA, see appendix I.)
2. PSU's should not cross regional lines, i.e., the four standard Census regions—Northeast, North Central, South, and West. However, it was not possible to follow this rule entirely as eight SMSA's crossed regional boundaries.
3. The area of a PSU should not exceed 2,000 square miles in the West Region and 1,500 square miles in other regions, except in cases where a single county exceeds the maximum area.
4. The 1960 population of each PSU should be at least 7,500 in the West Region and 10,000 in other regions, except in cases where this would require exceeding the specified maximum area.
5. PSU's should be formed in such a way as to avoid extreme length in any direction.
6. For situations in which more than one county was to be grouped to form a PSU, the principle was to make the groups as heterogeneous as possible with respect to a number of variables. The principal ones were economic area, principal industry (used primarily in urban areas), value of agricultural products (used primarily in rural areas), and the proportion of the county's population that was not white. The last item was used only in areas where there was appreciable variation between counties, primarily in the South.

A more detailed description of the formation of the PSU's may be found in Bureau of the Census' Technical Paper No. 7.<sup>21</sup>

#### Stratification of PSU's

The sample designs for CPS and HIS have changed several times since the surveys began, but in 1962 when Cycle II was designed both consisted of 357 strata and 357 sample PSU's—one PSU from each stratum. In determining which PSU's should be grouped together to form a stratum, a number of factors were considered.

1. Since only one PSU was to be selected from a stratum with a probability proportional to a measure of size, each PSU with a population above a certain size was put into a separate stratum by itself. Those PSU's are referred to as "self-representing." The population size cutoff for self-representing PSU's when most of the stratification work was done in the 1950's was 400,000 according to the 1950 population census. In some instances, however, a PSU with less than 400,000 people was classified as self-representing. These were smaller SMSA's within 100 miles of an SMSA with over 400,000 people. This was done since the field organization that served the larger city could also serve the smaller one and thus reduce survey costs.

In 1962 when the HIS sample was redesigned, utilizing 1960 census data, an additional criterion was introduced; namely, that PSU's with a population size greater than 75 percent of the national average for non-self-representing strata should also be self-representing. The result of this was that all PSU's with more than 242,000 population in 1960 were classed as self-representing. For the 357 strata, 112 PSU's were self-representing and 245 sample PSU's were not.

2. Strata should be approximately the same size except where a single PSU was larger than



an average stratum. The average population for non-self-representing strata within regions ranged from 298,000 to 349,000 (table I).

3. Strata containing more than one PSU would be as homogeneous as possible. Combining this with the principle for forming PSU's, a stratum should contain PSU's which tend to be alike, but the ultimate sampling units within PSU's should be as unlike as possible. The basic modes of stratification were:

- SMSA or not

- Rate of population change, 1950 to 1960

- Percent of population living in urban areas

- Percent of population in manufacturing

- Principal industries

- Average value of retail trade

- Proportion of population that was not white

4. The geographic spread of PSU's for non-self-representing strata is restricted only by the four census regional boundaries. That is, a stratum might be composed of PSU's located anywhere in a region but cannot contain non-self-representing PSU's located in different regions. Some effort was made, however, to combine PSU's located in the same Census division within regions.

The first step in the stratification process was to allocate each PSU to one of three groups. All self-representing PSU's were assigned to group 1; non-self-representing PSU's located in areas of relatively high population density were put in group 2, and the remaining PSU's were assigned to group 3. The next step was to classify groups 2 and 3 into three groups according to degree of urbanization. One subgroup contained SMSA's not classified in group 1. The other two subgroups were labeled "urban" and "rural." A PSU was considered rural if its rural farm population was 35 percent or more of the total, or if the rural farm population of the PSU was less than 35 percent but the population in urban places was less than the rural farm population and the rate of population increase was well below the average for the general area in which the PSU was located. After those two steps, stratification proceeded with primary attention being given to rate of population increase, degree of urbanization, color, principal industry, and type of farming. After semi-final stratification was completed the results were reviewed, and a few subjective changes were made which reviewers thought would increase socioeconomic homogeneity between PSU's within strata. Thus 357 strata were formed which have characteristics as shown in table I.

Table I. Number and average size of strata in the 357 area design by type of strata and region

Region and type of strata	Number of strata	Average 1960 strata population	SMSA		Non-SMSA	
			Number of strata	Total 1960 population <sup>1</sup> (in thousands)	Number of strata	Total 1960 population (in thousands)
<u>Self-representing</u>						
Total <sup>2</sup> -----	112	898,000	107	99,228	5	1,296
Northeast-----	28	1,225,000	23	32,905	5	1,296
North Central-----	26	1,013,000	26	26,346	-	-
South-----	36	571,000	36	20,563	-	-
West <sup>2</sup> -----	22	882,000	22	19,415	-	-
<u>Non-self-representing</u>						
Total <sup>2</sup> -----	245	322,000	39	13,515	206	65,283
Northeast-----	30	349,000	7	2,785	23	7,691
North Central-----	76	333,000	13	4,614	63	20,659
South-----	110	313,000	17	5,399	93	29,012
West <sup>2</sup> -----	29	298,000	2	717	27	7,921

<sup>1</sup>Because of minor differences between HIS design and Census in what was treated as an SMSA, the total of SMSA population on the table is about 141,000 less than SMSA total according to 1960 Census of Population.

<sup>2</sup>Includes Alaska and Hawaii.

— ○ ○ ○ —

APPENDIX III

HOUSEHOLD QUESTIONNAIRE

**CONFIDENTIAL** - The National Health Survey is authorized by Public Law 652 of the 84th Congress (70 Stat. 189; 42 U.S.C. 305). All information which would permit identification of the individual will be held strictly confidential, will be used only by persons engaged in and for the purposes of the survey and will not be disclosed or released to others for any other purposes (22 FR 1687).

**BUDGET BUREAU NO. 68-R620-S4.5**  
**APPROVAL EXPIRES JULY 31, 1965**

FORM NHS-HE5-2 (11-13-63) U.S. DEPARTMENT OF COMMERCE BUREAU OF THE CENSUS ACTING AS COLLECTING AGENT FOR THE U.S. PUBLIC HEALTH SERVICE

**NATIONAL HEALTH SURVEY**

1. Questionnaire of \_\_\_\_\_ Questionnaires

2. (a) Address or description of location (include city, zone, and State) \_\_\_\_\_

3. Identification code \_\_\_\_\_ 4. PSU number \_\_\_\_\_ 5. Segment number \_\_\_\_\_ 6. Serial number \_\_\_\_\_

2. (b) Mailing address if not shown in 2(a) OR  Same as shown in 2(a) \_\_\_\_\_

If this questionnaire is for an "EXTRA" unit in a B or NTA Segment, enter:

Serial No. of original Sample Unit \_\_\_\_\_ Item No. by which found \_\_\_\_\_ If in NTA Segment, also enter for FIRST unit listed on property \_\_\_\_\_

Segment List Sheet No. \_\_\_\_\_ Line No. \_\_\_\_\_

2. (c) Name of special dwelling place \_\_\_\_\_ Code \_\_\_\_\_ 7. Type of living quarters (Check one box)  Housing unit  Other unit

**L** Ask items 8 and 9 only if "Rural" box is marked  Rural  All other (Skip to Item 10)

8. Do you own or rent this place?  
 Own (Ask 9(a))  Rent (Ask 9(b))  Rent free (Ask 9(a))

9. (a) If Own or Rent free, ask - Does this place have 10 or more acres?  
 Yes  No

(b) If Rent, ask - Does the place you rent have 10 or more acres?  
 Yes  No

(c) During the past 12 months did sales of crops, livestock, and other farm products from the place amount to \$50 or more?  
 Yes  No

(d) During the past 12 months did sales of crops, livestock, and other farm products from the place amount to \$250 or more?  
 Yes  No

ALL segments (ask if Item 2(a) address identifies a SINGLE-UNIT structure).  
 10. Are there any occupied or vacant living quarters BESIDES YOUR OWN --  
 -- in the basement? ....  Yes--S \_\_\_\_\_ L \_\_\_\_\_  No  
 -- on this floor? .....  Yes--S \_\_\_\_\_ L \_\_\_\_\_  No  
 -- on any other floor of this building? ....  Yes--S \_\_\_\_\_ L \_\_\_\_\_  No  
 (Fill Table X for each quarters NOT listed)

ALL segments (ask if Item 2(a) identifies entire floor or unnumbered part of floor in a MULTI-UNIT structure).  
 11. Are there any occupied or vacant living quarters BESIDES YOUR OWN --  
 If Item 2(a) identifies entire floor -- on this floor?  Yes--S \_\_\_\_\_ L \_\_\_\_\_  No  
 If Item 2(a) identifies part of the floor, specify part -- in the -- of this floor? (Fill Table X for each quarters NOT listed.)

TA and NTA segments (ask at all units EXCEPT APARTMENT HOUSES).  
 12. Is there any other building on this property for people to live in - either occupied or vacant?  
 Yes--S \_\_\_\_\_ L \_\_\_\_\_  No  
 (Fill Table X for each quarters NOT listed.)

Telephone No. \_\_\_\_\_  
 13. What is the telephone number here? \_\_\_\_\_ OR  No telephone

(INTERVIEWER): If eligible child in household enter child's name, segment, serial, and column number on Medical History Form.  
 (READ TO RESPONDENT)  
 In addition to the information you have already given me, I would like to leave this form to be filled out about -- . The form is self-explanatory. A representative of the U.S. Public Health Service will come by to pick up the form in a week or so. (Ask Item 14)

14. What would be the best time of day for the representative to come?.....  
 Medical histories left for-- \_\_\_\_\_ Person with whom form left-- \_\_\_\_\_  
 Column No(s). \_\_\_\_\_ Column No. and relationship \_\_\_\_\_

15. RECORD OF CALLS AT HOUSEHOLD

Item	Date	1		2		3		4		5	
			Com.		Com.		Com.		Com.		Com.
Entire household	Time										

16. REASON FOR NON-INTERVIEW

TYPE	A	B	C	Z
Reason:	<input type="checkbox"/> Refusal (Describe in footnotes) <input type="checkbox"/> No one at home-- repeated calls <input type="checkbox"/> Temporarily absent (Go to 17) <input type="checkbox"/> Other (Specify)	<input type="checkbox"/> Vacant--non-seasonal <input type="checkbox"/> Vacant--seasonal <input type="checkbox"/> Usual residence elsewhere <input type="checkbox"/> Other (Specify)	<input type="checkbox"/> Demolished <input type="checkbox"/> In sample by mistake <input type="checkbox"/> Eliminated in sub-sample <input type="checkbox"/> Other (Specify)	Interview not obtained for Cols. _____ because:

17. TYPE A FOLLOW-UP PROCEDURE  
 Final call results in a Type A non-interview (except Refusals) take the following steps:  
 1. Contact neighbors (caretakers, etc.) until you find someone who knows the family.  
 2. Find out the number of people in the household, their names and approximate ages; if names of all members not known, ascertain relationships. Record this information in the regular spaces inside the questionnaire.

18. Signature of interviewer \_\_\_\_\_ 19. Code \_\_\_\_\_

ALL	1. (a) What is the name of the head of this household? (Enter name in first column.) (b) What are the names of all other persons who live here? (List all persons who live here.) (c) I have listed (Read names) is there anyone else staying here now such as friends, relatives, or roomers? ..... <input type="checkbox"/> Yes (List) <input type="checkbox"/> No (d) Have I missed anyone who usually lives here but is now -- Temporarily in a hospital? <input type="checkbox"/> Yes (List) <input type="checkbox"/> No -- Away on business? ..... <input type="checkbox"/> Yes (List) <input type="checkbox"/> No -- On a visit or vacation? ... <input type="checkbox"/> Yes (List) <input type="checkbox"/> No (e) Do any of the people in this household have a home anywhere else? <input type="checkbox"/> Yes (Apply household membership rules, if not a household member delete) <input type="checkbox"/> No (Leave on questionnaire)	Last name <b>1</b> ----- First name -----
	2. How are(is) -- related to the head of the household? (Enter relationship to head, for example: wife, daughter, stepson, grandson, mother-in-law, partner, roomer's wife, etc.)	Relationship <b>HEAD</b>
	3. Race (Mark one box for each person)	<input type="checkbox"/> White <input type="checkbox"/> Negro <input type="checkbox"/> Other
	4. Sex (Mark one box for each person)	<input type="checkbox"/> Male <input type="checkbox"/> Female
	5. (a) How old were you on your last birthday?  For each child age 5-12 listed on the questionnaire, ask: (b) What is the month, day, and year of --'s birth? (Check with Question 5(a) for consistency)	Age <input type="checkbox"/> Under 1 year Month Day Year
TO INTERVIEWER: Mark "EC" box for each eligible child (age 6-11) listed on the questionnaire. If no EC, ask coverage questions on Page 1. NOTE: Questions 6-14 must be asked only of parent(s) or guardian(s) of EC. If no parent or guardian is at home, arrange to call back when they will be home.		<input type="checkbox"/> EC <input type="checkbox"/> Not EC
ASK FOR EC	6. What is the name and location of the school -- goes to?  (a) What grade is -- in?	<input type="checkbox"/> No school Name and location  Grade
	7. Where were you born? (Check U.S. box or write in name of country)	<input type="checkbox"/> U.S. Foreign country
ASK FOR PARENTS OR GUARDIANS OF EC	8. Are you primarily right handed, primarily left handed, or both?	<input type="checkbox"/> Right <input type="checkbox"/> Left <input type="checkbox"/> Both
	9. What is the highest grade you attended in school? (Circle highest grade attended or mark "None.") (If attended, ask): (a) Did you finish this grade (year)?	<input type="checkbox"/> None Elem.... 1 2 3 4 5 6 7 8 High.... 1 2 3 4 College 1 2 3 4 5+ <input type="checkbox"/> Yes <input type="checkbox"/> No
	10. What were you doing most of the past 3 months -- working, keeping house, or doing something else? (If "Doing something else," ask): (a) What were you doing? (Enter reply verbatim and ask 10(b)) ..... (If "Keeping house" OR "Doing something else," ask): (b) Did you work at a job or business at any time during the past 3 months? ..... (If "Working" in 10 OR "Yes" in 10(b), ask): (c) Did you work full-time or part-time? .....	<input type="checkbox"/> Working <input type="checkbox"/> Keeping house <input type="checkbox"/> Something else ----- <input type="checkbox"/> Yes <input type="checkbox"/> No ----- <input type="checkbox"/> Full-time <input type="checkbox"/> Part-time
11. Are you now married, widowed, divorced, or separated? (If "Married," ask): (a) Have you (your husband) been married more than once?	<input type="checkbox"/> Married <input type="checkbox"/> Divorced <input type="checkbox"/> Widowed <input type="checkbox"/> Separated <input type="checkbox"/> Yes <input type="checkbox"/> No	
PARENTS ONLY	12. Besides (Read names of children entered in Question 1) have you and/or your husband(wife) ever had any other children? <input type="checkbox"/> Yes <input type="checkbox"/> No <input type="checkbox"/> No parent (If "Yes," ask): (a) What are their names? (b) How old is --? (If now deceased enter date of birth) (c) Where does he/she live now? (If now deceased enter "deceased")	Name ----- ----- -----
	13. Please look at this card (Hand respondent HES-2(a) card and pencil). Do any of the questions on that card apply to any members of the family? Please mark "Yes" or "No" for each question. (For each "Yes" marked, ask): (a) You have checked --. Who was this? (b) When was this? <div style="border: 1px solid black; padding: 2px; display: inline-block;">NOTE: If "1" marked, enter name of hospital or institution.</div>	Statement No. ----- ----- -----
ALL EC HOUSEHOLDS	14. Which of these income groups represents your total combined family income for the past 12 months, that is, your's, your --'s, etc? (Show Income Flash Card HES-2(b).) Include income from all sources, such as wages, salaries, rents from property, Social Security, or retirement benefits, help from relatives, etc. (Go to Question 15 on Page 4)	Group

Last name (2)			Last name (3)			Last name (4)			Last name (5)			Last name (6)		
First name			First name			First name			First name			First name		
Relationship			Relationship			Relationship			Relationship			Relationship		
<input type="checkbox"/> White <input type="checkbox"/> Negro <input type="checkbox"/> Other			<input type="checkbox"/> White <input type="checkbox"/> Negro <input type="checkbox"/> Other			<input type="checkbox"/> White <input type="checkbox"/> Negro <input type="checkbox"/> Other			<input type="checkbox"/> White <input type="checkbox"/> Negro <input type="checkbox"/> Other			<input type="checkbox"/> White <input type="checkbox"/> Negro <input type="checkbox"/> Other		
<input type="checkbox"/> Male <input type="checkbox"/> Female			<input type="checkbox"/> Male <input type="checkbox"/> Female			<input type="checkbox"/> Male <input type="checkbox"/> Female			<input type="checkbox"/> Male <input type="checkbox"/> Female			<input type="checkbox"/> Male <input type="checkbox"/> Female		
Age <input type="checkbox"/> Under 1 year			Age <input type="checkbox"/> Under 1 year			Age <input type="checkbox"/> Under 1 year			Age <input type="checkbox"/> Under 1 year			Age <input type="checkbox"/> Under 1 year		
Month	Day	Year	Month	Day	Year	Month	Day	Year	Month	Day	Year	Month	Day	Year
<input type="checkbox"/> FC <input type="checkbox"/> Not EC			<input type="checkbox"/> EC <input type="checkbox"/> Not EC			<input type="checkbox"/> EC <input type="checkbox"/> Not EC			<input type="checkbox"/> EC <input type="checkbox"/> Not EC			<input type="checkbox"/> EC <input type="checkbox"/> Not EC		
<input type="checkbox"/> No school			<input type="checkbox"/> No school			<input type="checkbox"/> No school			<input type="checkbox"/> No school			<input type="checkbox"/> No school		
Name and location			Name and location			Name and location			Name and location			Name and location		
Grade			Grade			Grade			Grade			Grade		
<input type="checkbox"/> U.S.			<input type="checkbox"/> U.S.			<input type="checkbox"/> U.S.			<input type="checkbox"/> U.S.			<input type="checkbox"/> U.S.		
Foreign country			Foreign country			Foreign country			Foreign country			Foreign country		
<input type="checkbox"/> Right <input type="checkbox"/> Left <input type="checkbox"/> Both			<input type="checkbox"/> Right <input type="checkbox"/> Left <input type="checkbox"/> Both			<input type="checkbox"/> Right <input type="checkbox"/> Left <input type="checkbox"/> Both			<input type="checkbox"/> Right <input type="checkbox"/> Left <input type="checkbox"/> Both			<input type="checkbox"/> Right <input type="checkbox"/> Left <input type="checkbox"/> Both		
<input type="checkbox"/> None Elem... 1 2 3 4 5 6 7 8 High... 1 2 3 4 College 1 2 3 4 5+			<input type="checkbox"/> None Elem... 1 2 3 4 5 6 7 8 High... 1 2 3 4 College 1 2 3 4 5+			<input type="checkbox"/> None Elem... 1 2 3 4 5 6 7 8 High... 1 2 3 4 College 1 2 3 4 5+			<input type="checkbox"/> None Elem... 1 2 3 4 5 6 7 8 High... 1 2 3 4 College 1 2 3 4 5+			<input type="checkbox"/> None Elem... 1 2 3 4 5 6 7 8 High... 1 2 3 4 College 1 2 3 4 5+		
<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No		
<input type="checkbox"/> Working <input type="checkbox"/> Keeping house <input type="checkbox"/> Something else			<input type="checkbox"/> Working <input type="checkbox"/> Keeping house <input type="checkbox"/> Something else			<input type="checkbox"/> Working <input type="checkbox"/> Keeping house <input type="checkbox"/> Something else			<input type="checkbox"/> Working <input type="checkbox"/> Keeping house <input type="checkbox"/> Something else			<input type="checkbox"/> Working <input type="checkbox"/> Keeping house <input type="checkbox"/> Something else		
<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Yes <input type="checkbox"/> No		
<input type="checkbox"/> Full-time <input type="checkbox"/> Part-time			<input type="checkbox"/> Full-time <input type="checkbox"/> Part-time			<input type="checkbox"/> Full-time <input type="checkbox"/> Part-time			<input type="checkbox"/> Full-time <input type="checkbox"/> Part-time			<input type="checkbox"/> Full-time <input type="checkbox"/> Part-time		
<input type="checkbox"/> Married <input type="checkbox"/> Divorced <input type="checkbox"/> Widowed <input type="checkbox"/> Separated <input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Married <input type="checkbox"/> Divorced <input type="checkbox"/> Widowed <input type="checkbox"/> Separated <input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Married <input type="checkbox"/> Divorced <input type="checkbox"/> Widowed <input type="checkbox"/> Separated <input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Married <input type="checkbox"/> Divorced <input type="checkbox"/> Widowed <input type="checkbox"/> Separated <input type="checkbox"/> Yes <input type="checkbox"/> No			<input type="checkbox"/> Married <input type="checkbox"/> Divorced <input type="checkbox"/> Widowed <input type="checkbox"/> Separated <input type="checkbox"/> Yes <input type="checkbox"/> No		
Age			Present whereabouts											
Name			Relationship			Year(s)			Name of Institution					
Group			Group			Group			Group					

15. Is any language other than English spoken here in your home?

Yes

No

(If "Yes," ask):

What language(s)?

Language(s) spoken \_\_\_\_\_

(Complete front page of questionnaire)

Comments: (Include here any information which might be useful to the PHS representative when she calls to pick up the Medical History Form.)

TABLE X - LIVING QUARTERS DETERMINATIONS AT LISTED ADDRESS

Line No.	Questionnaire No.	Are these (Specify location) quarters for more than one group of people?		Location of unit (Examples: Basement, 2nd floor, etc.)	USE OF CHARACTERISTICS						CLASSIFICATION		IF HU IN B SEGMENT, ASK		
		Yes (Fill one line for each group)	No		Occupied		All Quarters				Not a separate unit (Add occupants to this questionnaire)	Fill separate questionnaire and interview		In what year were these (Specify location) quarters created? (If 1959 or 1960, also specify "F" if first half or "L" if last half)	(If before July 1960) What was the name of the household head of these quarters on April 1, 1960?
					Do the occupants of these (Specify location) quarters live and eat with any other group of people?	Do these (Specify location) quarters have:	Direct access from the outside or through a common hall?	A kitchen or cooking equipment for exclusive use?	Yes	No		Yes	No		
(1)	(2)	(3a)	(3b)	(4)	Yes (5a)	No (5b)	Yes (6a)	No (6b)	Yes (7a)	No (7b)	(8)	(9a)	(9b)	(10)	(11)
1															
2															

## VITAL AND HEALTH STATISTICS PUBLICATION SERIES

*Formerly Public Health Service Publication No. 1000*

- Series 1. Programs and collection procedures.*—Reports which describe the general programs of the National Center for Health Statistics and its offices and divisions, data collection methods used, definitions, and other material necessary for understanding the data.
- Series 2. Data evaluation and methods research.*—Studies of new statistical methodology including: experimental tests of new survey methods, studies of vital statistics collection methods, new analytical techniques, objective evaluations of reliability of collected data, contributions to statistical theory.
- Series 3. Analytical studies.*—Reports presenting analytical or interpretive studies based on vital and health statistics, carrying the analysis further than the expository types of reports in the other series.
- Series 4. Documents and committee reports.*—Final reports of major committees concerned with vital and health statistics, and documents such as recommended model vital registration laws and revised birth and death certificates.
- Series 10. Data from the Health Interview Survey.*—Statistics on illness, accidental injuries, disability, use of hospital, medical, dental, and other services, and other health-related topics, based on data collected in a continuing national household interview survey.
- Series 11. Data from the Health Examination Survey.*—Data from direct examination, testing, and measurement of national samples of the civilian, noninstitutional population provide the basis for two types of reports: (1) estimates of the medically defined prevalence of specific diseases in the United States and the distributions of the population with respect to physical, physiological, and psychological characteristics; and (2) analysis of relationships among the various measurements without reference to an explicit finite universe of persons.
- Series 12. Data from the Institutional Population Surveys.*—Statistics relating to the health characteristics of persons in institutions, and their medical, nursing, and personal care received, based on national samples of establishments providing these services and samples of the residents or patients.
- Series 13. Data from the Hospital Discharge Survey.*—Statistics relating to discharged patients in short-stay hospitals, based on a sample of patient records in a national sample of hospitals.
- Series 14. Data on health resources: manpower and facilities.*—Statistics on the numbers, geographic distribution, and characteristics of health resources including physicians, dentists, nurses, other health occupations, hospitals, nursing homes, and outpatient facilities.
- Series 20. Data on mortality.*—Various statistics on mortality other than as included in regular annual or monthly reports—special analyses by cause of death, age, and other demographic variables, also geographic and time series analyses.
- Series 21. Data on natality, marriage, and divorce.*—Various statistics on natality, marriage, and divorce other than as included in regular annual or monthly reports—special analyses by demographic variables, also geographic and time series analyses, studies of fertility.
- Series 22. Data from the National Natality and Mortality Surveys.*—Statistics on characteristics of births and deaths not available from the vital records, based on sample surveys stemming from these records, including such topics as mortality by socioeconomic class, hospital experience in the last year of life, medical care during pregnancy, health insurance coverage, etc.

For a list of titles of reports published in these series, write to:

Office of Information  
National Center for Health Statistics  
Public Health Service, HSMHA  
Rockville, Md. 20852