Irving Sivin, New York City Health Department Paul M. Densen, New York City Health Department

#### I INTRODUCTION

In the first week of January, 1964, New York City's Department of Health began a probability sample survey of the city's residents. Its chief purpose was to provide hitherto unavailable data, on a continuing basis, about the health status of the population of the city. Because the Health Department is an important provider of medical services to the community, we wanted to obtain a clear picture of how the public and certain groups within it obtain their medical care. Consequently we focused on the area of medical economics. Our questions were designed to determine not only the amount of medical care received by New Yorkers and the diseases for which this care was required, but also to determine the auspices under which it is given, and how families finance their care.

The survey was conceived as a necessary supplement to the traditional vital and service statistical activities of the Department. We hoped that the data obtained in the survey would prove useful in measuring the effectiveness of the Department's current service programs, and could aid in planning new programs. We also felt that the creation of a sampling unit within the Department would provide it with the materials and skills required to conduct AD HOC surveys as the need arose.

The decision to undertake a continuing health survey of relatively large dimensions was also made in the belief that the data collected would prove useful to other city agencies as well. Reliable current data on population, household size, family income, migration, etc. are not available in late inter-Censal years on a local basis, although there is considerable need for them. Because the Population Health Survey would collect such data on a routine basis, we felt the responsibility to make these measurements with as great precision as possible.

It is the purpose of this paper to describe the sample design and data processing controls which were used to achieve the goals set for the survey.

II SAMPLE DESIGN

A. FRAMEWORK FOR THE DESIGN

Sample design should be married to the purposes for which a survey is undertaken, and should utilize the available resources in the most efficient manner. The initial decision to undertake a health survey of the City was made against a background of limited resources to defray the costs of interviewing, data processing and publication of the data. The survey staff would be regular Civil Service employees of the Department of Health. Detailed sample design had to fit into this framework.

It was decided that the most useful way to investigate the health status of the population in the context of a household survey was to ask questions similar to those used in the National Health Survey. The broad outline of the sample design was also to follow the National Health Survey's design. New York's survey would be an annual sample, divided into 52 equal subsamples. One of these subsamples would be interviewed each week throughout the year. From one year to the next, a different, but neighboring, set of households would be interviewed.

The population to be surveyed was the civilian non-institutional population of the City.

B. DETERMINATION OF SAMPLE SIZE

The size of our budget, when set against the expected cost of interviewing, indicated that the upper limit of our sample size would be about 7500 households per year. Our principal problem was to determine the minimum number of households that would satisfy our needs for reliable detailed data on an annual basis.

The economics of interviewing and the resources at our disposal demanded that we employ a cluster sample. Therefore we had to explore the relationships between the size of cluster and its sampling efficiency, and the costs of interviewing together with the cost of the delineation of clusters. We judged, on the basis of response data available to us from the Washington Heights Survey in New York, that the cost per interview would be relatively invariant for clusters of four or more households and would be about \$6 per household. We also concluded that the cost of the delineation of clusters would be proportional to the number of clusters in sample, about \$5 per cluster. We interpreted variance data from the National Health Survey to mean that for health characteristics, larger clusters would not seriously inflate the sampling errors, when compared to the cost advantages of having fewer of them. Our principal concern about large clusters lay in their inefficiency with regard to socio-economic information.

Before making a decision with regard to the size of the cluster to be employed, we decided to establish a list of the key statistics which were to be derived from the survey. We hoped that this process would indicate, in terms of simple random sampling, how large a sample we should have, and from this vantage we could extrapolate to the size of cluster sample (and of cluster) that we needed.

The list of the survey's chief concerns was based on questions submitted to us by the heads of major units within the Health Department. Our survey sought to answer the following questions:

- 1. How do different income and ethnic groups within the City finance their medical care? For these groups, what proportion of physicians' services are financed out of pocket?
- 2. What proportion of families have hospital insurance for all members of the family? How does this proportion vary by family size, ethnic group and income?

- 3. What proportion of families have one or more family members hospitalized during the course of the year? How is the length of stay affected by family size, income, hospital insurance and type of disease?
- 4, What proportion of out-patient medical services in the City are provided by governmental agencies?
- 5. How many physically handicapped persons are there in New York?
- 6. How many physician visits are made per person per year? How many dental visits are made?

These questions implied that comparisons would be made between the characteristics of different ethnic and socio-economic groups. Our goal therefore was to be reasonably certain that valid comparisons could be made. We felt that this goal would be achieved if an observed difference of 10% in a characteristic for two different socio-economic groups would prove statistically significant.

The smallest ethnic group for which the Department wished detailed information was the Puerto Rican population, which constituted about 8% of the City's population in 1960. The next smallest ethnic group for which detailed data were desired was the Negro population. There are about twice as many Negroes as there are Puerto Ricans in New York. A simple random sample of families would yield about twice as many Negro families in sample as Puerto Rican families, and the reliability conditions outlined above would be satisfied approximately by the equation

$$\frac{PQ}{P} + \frac{PQ}{2n} = (.05)^2$$

Since the highest value of PQ is .25, a simple random sample that yielded 150 Puerto Rican families and 300 Negro families would suffice. The total sample size needed to produce these numbers would be 1875 families. Were cluster sampling to prove only one-third as efficient as simple random sampling in obtaining the economic characteristics of these groups, then a sample of about 5625 families would be required. We speculated that cluster sampling would operate at about this efficienty, if the cluster size were between 6 and 10 units.

A somewhat different orientation was given to the problem of sample size when we realized that data from the survey would be used as an adjunct, if not the principal source, in obtaining current estimates of the total population of the City. We felt that estimates from the survey would prove useful if the total population could be estimated with a relative error of 2% or less, that is, with a two sigma error of about 300,000 persons. We felt this to be adequate because by the end of 1963 the difference between the Census Bureau's projections of the population and those of the City Planning Commission had risen to about 280,000 persons.

To achieve a relative error of 2% on the estimate of the total, as well as for the sake of providing uniform interviewer assignments, it was necessary to make the clusters contain approximately equal numbers of housing units and of

persons. Doing this would diminish the variance between cluster totals which is the principal source of sampling error in estimates of the population. The first and chief step in controlling the size of cluster could be made by using Census Blocks as our primary sampling units and selecting these with unequal probabilities, proportional to the number of units reported in them by the Census in 1960.

To achieve a relative error of 2% for the population total, we felt that it would be necessary to achieve an error of about 1% on the estimate of the total number of housing units in the City. If we could obtain a coefficient of variation of the size of cluster of .25 through the use of PPS sampling and careful delineation of clusters, then 625 clusters would suffice to yield an error of 1% for the entire sample. We felt that it would be possible to do this using clusters of about eight units. It seemed to us that were the cluster size smaller, there would be a higher underlying coefficient of variation, which would necessitate the inclusion of more clusters in sample.

We, of course, had to take into account not only the variation of the number of units per cluster, but also the variation of the number of people within the units. We assumed that the relative variance of the estimated population total would be a function of the form



where  $V^{\mathbf{2}}$  is the relative variance

is the within-cluster correlation co-efficient of persons per unit.

efficient of persons per unit.

- n is the average number of units per
- cluster.

We assumed that the relative variance of the number of persons per housing unit was equal to 0.4, a value somewhat higher than the figure for the number of persons per occupied unit, which we derived from the 1960 Census. Under the assumption that cluster sampling would be only 50% as efficient for this characteristic as simple random sampling,  $1 + \int (\pi - 1)$  equals 2.0. Therefore we would expect that 625 clusters of 8 units would yield a relative error of 1.6% on the estimate of the total population. Under the assumption that the efficiency would prove to be only 33% of simple random sampling, that is with  $1 + \int (\overline{n} - 1)$ equal to 3.0, a sample of 625 clusters of eight households would yield a relative error of 1.9%, on the estimated total population.

Ultimately, as a compromise between the indicated values, and for the sake of simplicity in estimation, we agreed on a sample with an overall fraction of one in 500. This would produce about 5700 housing units and 5400 households in sample each year. The sample would have about 700 clusters of eight housing units each.

### C. GEOGRAPHIC STRATIFICATION

We arrived at our sample size under a set of assumptions which are relevant to a simple random sample of clusters. Actually, from the very beginning we had envisaged that the Population Health Survey would be a stratified sample, in which the strata were to be the 30 Health Districts into which New York City is divided. Stratification would ensure that the representation of the various socio-economic and demographic groups in the City would be close to their level in the population. Because the Health Districts differ greatly in their demographic and socio-economic composition, they also differ greatly in the type and amount of public health services provided by the department. We desired to produce data, from time to time, on the health and medical care characteristics of groups heavily and lightly serviced districts, and therefore needed stratification.

The introduction of geographic stratification would not greatly affect the overall sample size we required for inter-group comparisons. For city-wide estimates, stratification along geographic lines, would produce at best only modest gains in the overall efficiency of the survey, and we felt no need to reduce the sample size on this score.

D. EFFECT OF ESTMATION PROCEDURES ON THE SURVEY DESIGN

The processing resources at our disposal precluded the use of all but the simplest estimation procedures for the data from the survey. Both for the sake of the efficiency of the sample and for simplicity in processing, we determined to have a self-weighting sample. All estimated totals derived from the survey would therefore, be the sample totals multiplied by the reciprocal of the overall sampling fraction. All estimated rates, percentages or proportions would have sample totals for both the numerator and the denominator. Since these simple techniques of estimation were to be employed, we could not count on any gains in the efficiency of the survey through their use.

E. NEW CONSTRUCTION STRATA

In a survey where the estimation procedures are simple inflations of survey results or just sample proportions, the actual efficiency achieved by the sample rests heavily on the amount of detailed work entering into the final design, and the degree to which the physical realization of the sample conforms to the blueprint.

To insure that the totals estimated from the sample have the precision we required, we realized that we had to create a separate selection frame, unrelated to the Census Block Statistics from which the sample in our geographic strata would be chosen. We needed a second sampling frame because New York adds about 30,000 new housing units to its inventory each year, a rate better than 1% per annum. Most of these additions are in the form of large apartment buildings or large developments of private homes. Were only one of these large developments to appear unexpectedly within a sample cluster, the precision of our estimate of the total number of housing units in the City would easily be cut in half. For example, let us assume that, aside from the presence of a single cluster containing 60 housing units, the sample would have achieved its goal: a coefficient of variation of the number of units per cluster of .25 on an individual cluster basis. Then, the presence of the single large cluster would have increased the relative variance of the sample of 700 clusters from .000089 to .000174, an increase of 94%. The argument is detailed below.

- Let  $\overline{\mathbf{x}} = 8$ . This is the average number of units per sample cluster, provided no large cluster is encountered.
- Let  $\mathbf{0}^2$  = 4. This is the variance of the number of units per cluster under the assumption that we have achieved a coefficient of variation of .25 in the size of the cluster. The relative variance is .0625 on an individual cluster basis.

The relative variance of a sample of 700 clusters under these assumptions is .000089 Now in a sample of 700 clusters which had achieved the above characteristics, the value of the sum of squares is 47,600, be-

$$\sum_{x_{i}^{2}/700}^{700} - 8^{2} = 4$$

cause

where  $X_i$  is the number of units in a cluster If one cluster with 60 units had been encountered instead of a cluster with 8 units, then the sum of squares in the sample would have increased by 3536 to 51,136.

The sample mean number of units would now be 8.08, not 8.00

The variance between the sample cluster values would increase and now be 7.93, not 4.00

The relative variance on an individual cluster basis would now be .1217, not .0625 The relative variance of the sample total would now be .0001738, not .000089.

The inflation of sampling error caused by a large cluster occurs not only in one or two characteristics, but is quite general. This is a result of the fact that the uniform rental or price structure of a development attracts a relatively homogenous population to it.

Since field work for the survey was to start almost four years after the Census date, we had to expect that we would encounter not one, but several clusters of new construction in the course of our work.

It was possible for us to create a separate selection frame for units built subsequent to the 1960 Census through the use of certificates of occupancy issued by the N.Y.C. Department of Buildings. These certificates indicate the address and the number of housing units contained in each new residential structure. They also contain a tax block number, which enabled us to determine how many new units were built in each tax block. Within each borough tax blocks were then selected for inclusion in the new construction sample with a probability proportionate to the number of new units in the block.

Since the sample derived from the 1960 Census is an area sample, it is necessary to prevent any new units which lie within the area clusters from being given a double chance of selection, should they be also represented in the new construction strata. Therefore each Census Block in the area sample was also identified by tax block number. These numbers are screened against the list of tax numbers in the new construction strata. New structures which appear in both the sample frames are then excluded from the area sample.

Every six months the new construction sample frame is updated. Since the Population Health Survey is a continuing survey, some of the newer units may have already been interviewed as part of an area cluster. If that has happened, the particular structure is not included in the updated new construction sample frame. This process assures us that the two sampling frames remain unduplicated. Clusters are selected from the updated frame with the requisite conditional probabilities, and are interviewed in the second half of the year.

F. ASSIGNMENT OF MEASURES TO CENSUS BLOCKS Since Census Blocks in New York City vary from zero to nine thousand units each, we could neither achieve an efficient sample nor maintain uniform interviewer workloads by giving each block an equal chance of coming into the sample. Therefore, we assigned to each block a number of measures, proportional to one-eighth of the Census count of housing units. Additional measures were assigned to those blocks which contained non-institutional group quarters, in order to keep the number of persons per measure relatively constant from block to block. Census Blocks with only a few reported units were amalgamated on the sample frame with adjacent blocks to prevent clusters with very few units from coming into sample. We also took pains to be sure that blocks which had not been recorded in Census tabulations (because they had no population in 1960) were given non-zero probabilities of selection. This was done by joining them to neighboring blocks.

The assignment of measures to individual blocks and the inclusion of all blocks in the sample frame were completely verified. The total number of measures assigned to blocks within a health district was also checked against a control number of measures, based on health district tabulations derived from the Census.

# G. SAMPLE SELECTION

In a stratified sample with a uniform sampling fraction, sample selection is usually made independently from stratum to stratum, and a simple random sample is taken within each stratum. We wanted to plan for the year to year change in our sample, for expansion of the sample in particular strata, and for the selection of special samples in connection with other Health Department studies. We also wanted to avoid, to as great a degree as possible, any tail-end variances which independent selection in our thirty odd strata might induce in the estimation of borough-wide and city-wide totals. A systematic sample of measures carried over from stratum to stratum was the best answer we could devise. Systematic sampling, however, has two inherent disadvantages. Periodicity in the population might be a multiple of the sampling interval, and therefore, a single systematic sample might be imprecise. Secondly, with systematic sampling, there is no unbiased way to estimate the sampling error. To overcome these difficulties we accumulated the block measures throughout the City in an order determined by three stages of randomization. In the first stage, a random permutation determined the order in which the five boroughs would have their measures cumulated. Within each borough another random permutation designated the order in which the health districts would appear in the cumulation. Finally, a third set of random permutations determined the order in which 357 small geographic areas would appear within the 30 health districts. Given this ordering, we employed just one random number and a systematic interval of 500 to select our sample of clusters for the entire City. Therefore the systematic effects, if any, could occur only within our smallest unit of randomization. We would still be left with a great number of degrees of freedom to estimate the sampling error because of the randomizing procedure. The whole process of sample selection in which we engaged may therefore be viewed as a single stage sampling. The cluster included in sample within a sample block corresponds to a translation of the cumulative random number which had selected the block.

Once the entire sample was selected for the initial survey year, weekly subsamples were established. We employed constraints in the subsampling process in order to insure that the geographic distance between the clusters interviewed in any two successive weeks would be as great as possible. As stated earlier, we scheduled a systematic half of the housing units in each cluster for interview during the first half of the year, and the other half for interview twenty-six weeks later.

## III REALIZATION OF THE SAMPLE A. CREATION OF CLUSTERS

Following the selection of a block, the address, inclusive of apartment number, of every residential unit in it was listed by our field staff. The count of the number of residential units was compared to the Census report, and the list of units was accepted when the count was within five percent of the Census value. Listings in disagreement with the Census report were reconciled by reference to Sanborn maps and to other information. Blocks with faulty listings or with irreconcilable ones were independently relisted, then reconciled.

Clusters were created out of the block listings so that the maximum size of any cluster within a block exceeded the minimum size by no more than three units. Wherever possible, we delineated compact clusters. In apartment houses to avoid ambiguity with regard to cluster boundaries and to keep the number of units per cluster constant, the ultimate sampling unit was often a systematic portion of a larger cluster.

The tightness of our control over the size of clusters was the final step which we could take to obtain the smallest sampling error from our survey. In the actual conduct of the survey all our efforts were bent to achieve the smallest possible biases in the data. We shall now turn to some of the efforts which we have put forward to control the mean square error of the Population Health Survey.

B. COVERAGE CHECKS

One of the important contributors to survey error is under-coverage of the target population. There are two components to undercoverage, missed households and missed persons within households. We have not designed a procedure to cope with the latter problem, but we do check on a sample basis for missed units within the clusters. To date, however, the best means of detecting the existence of missed units has been the second set of interviews taken in the sample clusters 26 weeks after the first set by a different interviewer. Published housing unit totals for New York City are in close agreement to our inflated sample values. We do not seem, therefore, to be suffering greatly from the under-coverage of units.

C. RESPONSE RATE

Perhaps the greatest contribution to the mean square error of the Population Health Survey has been made by non-response of sample households. In the first year of survey operations the response rate has proved to be 88% of the eligible households. In Manhattan, response has been 83%. While we are not chagrined by these results, they are well below the standard set by Census Bureau operations in New York. We are making every effort to improve our record in the second year of the survey operation.

D. INTERVIEWER EFFECT ON THE MEAN SQUARE ERROR

In terms of its weekly sample the Population Health Survey is a small operation, requiring only a few interviewers. To increase the number of people who collect the data, and thereby decrease the effect of a single interviewer on the mean square error on the statistics published, we have deliberately kept the weekly workload of each interviewer as small as possible. Because it has not proved feasible to assign interviewers to work outside their home boroughs, estimates of borough values represent the product of very few hands, and are subject to high variability on that account. Within the boroughs, however, we randomly assign clusters to the interviewers, thereby reducing the effect of interviewer-area interactions.

Our budget does not permit us to engage in a large re-interview program. It is therefore impossible for us to assess accurately the impact of so few interviewers on the data produced. Nevertheless, we find it useful to produce tabulations of health data by interviewer each quarter to discern gross differences in performance, particularly within boroughs where we have randomized assignments.

E. DATA PROCESSING CONTROLS

After the initial input of data onto the questionnaires, an inflation of the mean square error of the survey is bound to occur during the course of data processing. By setting up stringent quality control checks we have attempted to minimize this inflation.

The completed questionnaires are given routine check-in edits, and the fact that the interview has indeed taken place is verified, primarily by phone. The segment lists are checked to determine that all units scheduled for interview have been properly included in the sample.

Since the questionnaire employed in the survey is not an instrument which can readily be key punched, the information contained in it is transcribed, after coding, onto forms suitable for further processing. The professional staff of the survey reviews the medical and occupational coding on a sample basis. The accuracy of the transcription is also reviewed on a sample basis.

Once the data has been key punched onto cards, a computer program is employed to detect inconsistencies in the data. Inconsistencies within individual cards, inconsistencies between the cards of a single person, and inconsistencies between the cards for different persons within the same family are detected by this program. Before any tabulations are produced, all the errors detected by this program are corrected, and the program is re-run, to be sure that the indicated changes have been made. To date only four percent of all the cards processed have had detectable errors.

IV. ESTIMATION AND SAMPLING ERROR

A. ESTIMATION

As stated earlier (Section II D) we use only simple inflation estimates for totals. Sample values of proportions, rates or percentages are used directly for population values of the same types.

Two different inflation factors must be employed. The factors correspond to the two different sampling fractions used in collecting the data. In order to estimate statistics such as the number of persons in New York or the number who have been discharged from a hospital the sample data is simply multiplied by 500. To estimate population values of other types of data which are collected with a two week reference period, such as the number of physician visits or the number of dental visits made in the City in a year, the sample total must be multiplied by 13,000. This factor is the product of the basic inflation factor, 500 times 26, since the year is conceived as 26 two week periods.

The inflation factors are applied to the sample data only after adjustment has been made for complete non-response of eligible households. This adjustment is made by duplicating the information collected for a responding household in the same cluster as that which had the non-response. When no information is available about the characteristics of the non-responding household we have duplicated at random one of the interviewed households in the cluster. We have frequently obtained information about the number of persons in the non-responding household, and when this information is available we have selected for duplication a household within the same cluster that contains the same number of persons.

Ve have examined the results of this duplication process by tabulating health data with the duplicated cards in the deck, and by tabulating the same data without using the duplicated cards. Differences in such statistics as the percent of persons currently medically attended, the percent hospitalized in the past year, and the percent with a current limitation of activity have been trivial. We expected to find, as we did, that there were significant differences between the two sets of tabulations in economic data, since high income clusters in Manhattan's East Side have proved difficult to interview.

The failure of respondents to answer certain questions during the course of the interview, either because they do not know the answers or because, as in the case of income questions, they refuse to state an answer, increases the means square error of the survey and causes additional problems in estimation. When nonresponse is large - with income it has run to 8% of our respondents - we exhibit characteristics for the class of persons with income and the characteristic is used frequently in cross tabulations we set up a separate category for income unknown. However, when we seek to establish median family income figures, non-respondents are allocated to age - race - occupation groupings before estimation begins.

For items such as hospital insurance coverage the non-response rate has been under 1%, and we have simply considered the percentage of persons covered to be that number giving positive responses to our questions divided by the total number of persons in sample. For hospitalization rates, however, we have excluded persons with hospital status unknown before computing the percentage of persons hospitalized.

B. SAMPLING VARIANCES

Each statistic produced by the Population Health Survey has a specific sampling error. It is quite beyond our means to produce exact estimates of the error of each survey statistic. We have therefore limited ourselves to estimating the variances of seventy important totals and 100 key rates or percentages. From these variance tabulations we attempt to generalize our findings so that they serve as guides to the sampling errors of other items published by the survey.

The relative variance of an estimated total,  $\dot{x}$  derived from the survey is computed as  $V_{\dot{x}}^2 = .998 \sum_{x_{1,c}} \frac{A_{1,c} \dot{x}_{X_{1,c}} - (\dot{x}_{X_{1,c}})^2}{(\dot{x}_{1,c} \dot{x}_{X_{1,c}})^2} ((\dot{x}_{1,c} \dot{x}_{X_{1,c}})^2)^2$ Where x is the enumerated total value Where x

- hi for all elements in cluster i of stratum h.
  - K is the number of clusters in
  - h stratum h.
  - is the total number of strata, 34 of which 30 are geographic and 4 are new construction strata. .998 is the finite sampling correction factor.

For a ratio of the form  $\sum x_{1,i} / \sum y_{1,i}$ , the estimated relative variance is computed as  $V_{1}^{2} = V_{2}^{2} + V_{2}^{2} - 2 V_{1}^{2} G$ This form of estimation probably over-

states the true sampling error within the health districts (strata), because it treats the sample within the districts as a simple random sample of clusters. It is a measure of the variance between as well as within the smaller geographic areas used in the third stage of the randomizing process. We do not feel, however, that the overstatement would be large or important with regard to health characteristics.

The items for which we computed relative variances were chosen to represent families of items which we believed would have different sampling efficiencies. These families comprised data on demographic, economic, two-week condition, hospitalization and medical attendance data.

For each family of data we fit a curve of the form\_2

٧ź a + b/X=

X is an estimated total; a and b are values determined by minimizing the squared relative residuals of the function. Two or three iterations of the process are often necessary to produce a good fit. In using a curve of this form as in many other instances, we have followed the lead of the National Health Survey. V. EVALUATION OF THE SURVEY DESIGN FROM SURVEY

RESULTS

At this date, tabulations are available from the first six months of data collection in 1964. From the results it would seem that the survey is functioning at a level of efficiency somewhat higher than we had expected. The relative variances of some principal demographic statistics are shown below:

Item	Relative	Relative
	Variance	Error
Occupied Housing Units	.000149	1.2%
Family Heads	.000302	1.7%
All Persons	.000298	1.7%
Unrelated Individuals	.002370	4.9%
Non-White Family Heads	.005083	7.1%
17 11 01	• •	

From these figures it seems clear that we have achieved our goal of estimating the population of the City with a relative error of under 2%, while employing only half the number of interviews we had regarded as necessary. The chief source of our over-estimate of our requirements in this regard lay in our assumption about the homogeneity of the size of households within clusters.

On the other hand with regard to the health and medical care characteristics of small groups in the population, the sample size is none too large. For example in the first published report of the Population Health Survey the hospital insurance coverage of non-white persons in New York was estimated to be 50.7%, an estimate which was not significantly different from the 42.2% coverage reported for persons of Puerto Rican birth or parentage. These values are from data collected in the first six months of the Survey's operation, and differences as large as 8.5% should prove significant at the conclusion of a complete cycle of enumeration. Smaller differences near the 50% level between these two groups will not be significant.

### VI POSSIBLE REDESIGN OF THE SURVEY

In order to have more detailed information about certain groups in the City, it is possible to redesign the Survey somewhat, when the first two years of survey operation have produced information about the characteristics of the City as a whole. Under consideration is a plan to expand the sample in the six or seven health districts with the most severe public health problems. This expansion could be accomplished without changing our budget greatly, by reducing the sample size in the rest of the City to two-thirds of its present level. Such a design would still permit us to produce reliable city-wide estimates while sharpening our knowledge of the health status and medical care economics of the population in these districts. SUMMARY STATEMENT

The Population Health Survey has produced detailed health and demographic data about the City of New York which is not available from any other source. While accomplishing this, it has been able to design and realize samples for special studies both within the Health Department and for other governmental agencies of the City of magnitudes ranging up to 25,000 households. The creation of such an instrument for a local health agency in a community in which no other central data collection center exists is well worth its relatively small cost.