# ㊷ Allocation in stratified sampling

## By W. Edwards Deming

*How the problem arises.* Thus far we have learned how to compute the variance of a sum or mean under certain sampling conditions; and emphasis has been place on the importance of either attaining a prescribed variance at minimum cost, or, alternatively, of attaining the smallest possible variance for a given allowable budget.

However, so far no indication has been given concerning the most efficient distribution of a sampling amongest the several classes into which a universe may have been divided, in order to achieve such desirable aims. It would be possible by cut and try methods to increase the size of sampling in one class and decrease it in another, each time recomputing the total cost and expected variance, and thus eventually to arrive at a satisfactory allocation of the sample amongest the classes, but this is tedious and it would be much better to have an indication before hand of just what the desired aims depend on, and how to achieve them directly. The solution of this problem, the proper allocation amongest the several classes, is the aim.

(286)

<u>Formulation of the problem</u>. The symbol $\Sigma$ will be used here to denote the sum of some particular characteristic in the universe.

For example, $\Sigma$ might be the total inventory of nails or flour, the total annual sales of nails or flour or of all commodities combined; or it might be the total unemployed, or employed at so many hours per week, or the number of males or females or both under 18 and in school.

The universe of elements to be sampled might consist of a list of all the dealers of some specific or general type (all hardware dealers, all establishments selling flour at retail) within a given region. Or it might consist of a list of all the dwelling units, or of all the small areas into which the region could be divided for purposes of sampling. The dealers, dwelling units, or areas, are supposed to have been classified or stratified possibly on the basis of location and likely also on some other significant characteristic such as inventory, sales, or number of employees or number of occupied dwelling units reported at the last census. In each class at the time the sample is to be taken there will be an (unknown) average inventory or sales per dealer or per area; or an average number of people employed, unemployed, or in school, either per household or per area. If the Greek letter $\mu 1, \mu \cdots$ denote

these actual averages today, supposedly theoretically obtainable if a complete count were taken, then

$$\mathcal{E} = N_1 \mu_1 + N_2 \mu_2 + \cdots \text{ through all classes} \qquad (1)$$

wherein $N_1$ is the actual (supposed known) number of number in Class 1. But, of course, we can only provide estimates of the Greek symbols, and if $x$, $m_1$, $m_2$, $\cdots$ denote the sample estimates of $\mathcal{E}$, $\mu_1$, $\mu_2$, $\cdots$, the equation of estimate to be used here is

$$X = N_1 m_1 + N_2 m_2 + \cdots \qquad (2)$$

$m_1$, $m_2$, $\cdots$ are here the average inventories, sales, employed, etc. obtained from the samples in the various classes. These estimates are subject to sampling errors; in fact if $m_1$ is obtained by averaging the inventories of $n_1$ sample dealers drawn at random from class 1, then

$$\operatorname{Var} m_1 = \frac{N_1 - N}{N_1 - 1} \frac{\sigma_1^2}{n_1} \qquad \left[ P. \qquad \right] \quad (3)$$

There will be a like contribution of variance from each class, wherefore

$$\operatorname{Var} X = N_1^2 \operatorname{Var} m_1 + N_2^2 \operatorname{Var} m_2 + \cdots$$

$$= N_1^2 \frac{N_1 - n_1}{N_1 - 1} \frac{\sigma_1^2}{n_1} + N_2^2 \frac{N_2 - n_2}{N_2 - 1} \frac{\sigma_2^2}{n_2} + \cdots \quad (4)$$

It should be noted that the problem to
be solved here is associated with a par-
-ticular kind of sampling and estimation
, a random sample of $n_1$ dealers, house-
-holds, or areas is to be taken from Class
1, $n_2$ from Class 2, etc. Moreover, and
very importance, the estimate × in to
be formed first by straight averaging
to get $m_1$, $m_2$, ...,
followed then by multiplication by

$N_1$, $\mu_0$, $\mu_2$, ... as Eq. 2 indicates.
In symbols, if $x_{ij}$ denote the inventary
of Dealer $j$ in Class $i$, then

$$m_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij}$$  [ The $n_i$ dealers or dwelling
units are supposedly   ⑤
numbered $j = 1, 2, ..., n_i$
after the sample is drawn]

In practice, instead of forming these
average one by one, it is usuall more
expedition to form the term of Eq. 2
directly by noting that

$$N_i \, m_i = \frac{N_i}{n_i} \sum_{j=1}^{n_i} x_{ij} \qquad \left[ \text{Note that } \sum_{j=1}^{n_i} x_{ij} \text{ is} \right.$$

merely the sum of the (b)
inventories of the
sample dealers in Class $i$

Here it is the sampling interval $N_i / n_i$, or the
ratio of universe to sample size in Class $i$, that
is used as a multiplier of the sample inventory

$$\sum_{j=1}^{n_i} x_{ij} \, .$$

If some other procedure of estimation were introduc-
ed, for example, the ratio of the inventories of the
sample-dealers of Class I at two different dates
(the Census date and the sample date), then Eqs.
2 and 3 would be replaced by corresponding
but quite different equations, and all that
followed would be affected. That is to say,
optimum ~~alloc~~ allocation under one procedure of
sampling and estimating will ~~be~~ not be
optimum under another procedure. Thus, although
in this chapter we shall solve a very
important problem in optimum allocation,
~~under one~~ the solution should be regarded as
applicable only in the specific procedure descr-
-ibed. It should be said, though, that most
of the remarks made here will hold good in other
~~procedures~~ as well. (290)

A more complicated problem, first sampling areal units and then subsampling dwelling units from them, is treated in a book by Hansen, Hurwitz, and Deming entitled a chapter in *Population Sampling*, Bureau of the Census, 1947.

The cost of conducting the survey will be

$$C = n_1 c_1 + n_2 c_2 + \cdots$$

wherein $c_i$ is the cost of a questionnaire in Class $i$ (including collection, follow-up, and tabulation, with a pro-rated share of overhead).

The problems to be solved can be stated as follows:

Problem I. Find what sampling intervals $N_1/n_1$, $N_2/n_2$, etc. will produce minimum Var $X$ for a fixed allowable budget $C$.

Problem II. Find what sampling intervals $N_1/n_1$, $N_2/n_2$, etc. will produce a prescribed Var $X$ at minimum cost

Fortunatly, the two problems almost completely overlap, and we shall set off to solve the first

Solution of the first problem; cost prescribed. The solution, when finally arrived at, will determine a set of values of $n_1$, $n_2$, --- which constitute the best sizes of sample in the several strata (i.e., best in the sense of giving minimum Var $X$ for a given budget $C$). These will be termed the "equilibrium" values from the analogy with problems of equilibrium in mechanics

Let $\delta$ denote a variation from an equilibrium value. Then if $\mathrm{Var}\, x$ is at a proper minimum, its variation will be zero to within higher power of $\delta n_j$: i.e.,

$$\delta\, \mathrm{Var}\, x = \frac{\partial\, \mathrm{Var}\, x}{\partial n_1}\, \partial n_1 + \frac{\partial\, \mathrm{Var}\, x}{\partial n_2}\, \partial n_2 + \cdots = 0 \quad (8)$$

~~It will be observed that the $n_2$ procedure.~~

It will be observed that the procedure is to perform the differentiations exactly as if $n_1, n_2 \cdots$ were continuous variables, when in fact they can proceed only by steps of unity. This procedure can ~~proceed only~~ be relied upon to give results for $n_1\, n_2, \cdots$ correct within a unit; in fact the results (Eq. 13 ahead) are probably exact

Now differentiation of the first terms on the right-hand side of Eq. 4. gives

$$\frac{\partial}{\partial n_1}\, N_1^2\, \frac{N_1 - n_1}{N_1 - 1}\, \frac{\sigma_1^2}{n_1} = \frac{\partial}{\partial n_1}\, \frac{N_1^3\, \sigma_1^2}{(N_1 - 1)\, n_1}$$

$$= \frac{\partial}{\partial n_1}\, \frac{N_1^2\, \sigma_1^2}{N_1 - 1} = -\frac{N_1^3\, \sigma_1^2}{(N_1 - 1)\, n_1^2} - 0 \quad (9)$$

The second term on the right is 0 because $n_1$ is not contained in the term directly above it.

Every term on the right of Eq. 4 will contribute a derivative like the one just written; hence from Eq. 8 it follows that

$$\frac{N_2^3 \, \sigma_1^2}{(N_2-1)\,n_1^2} \, \delta n_1 + \frac{N_2^3 \, \sigma_2^2}{(N_2-1)\,n_2^2} \, \delta n_2 + \cdots = 0 \quad (10)$$

Now differentiate Eq. 7 to get

$$c_1 \, \delta n_1 + c_2 \, \delta n_2 + \cdots = 0 \quad ----- (11)$$

The $C$ on the right expresses the fact that however the sample sizes are permitted to vary, the total budget $C$ is to be kept fixed.

This equation thus imposes a condition on the variations $\delta n_1$, $\delta n_2$, ---; they can not all be independent; one of them is fixed by Eq. 11

Next, multiply Eq. 11 through by $-\lambda^2$ and add it to Eq. 10, collecting the coefficients of $\delta n_1, \delta n_2, ---$ The result is

$$\left[\frac{N_2^3 \, \sigma_2^2}{(N_2-1)\,n_1^2} - \lambda^2 c_2\right] \delta n_2 + \left[\frac{N_2^3 \, \sigma_2^2}{(N_2-1)\,n_2^2} - \lambda^2 c_2\right]$$

$$\delta n_2 + \cdots = 0 \quad ----- (12)$$

Now suppose I let $\delta n_2, \delta n_3 \ldots$ take on any values whatever, forcing Eq. 11 to hold by a proper choice of $\delta n_1$, at the same

time choosing $\lambda$ so as to make the coefficient of $\delta n_1$ vanish in Eq. 12; the first term is then 0 no matter what $\delta n_1$ may be. But this requires that the sum of all the remaining terms be 0 for any values whatever of $\delta n_2$, $\delta n_3$, ---, which can be so only if each coefficient is separately equal to 0. And thus every coefficient vanishes, giving

$$\lambda^2 c_i = \frac{N_i^2}{n_i^2}\,\sigma_i^2 \quad \frac{N_i}{N_i - 1} \quad \text{for all } i \text{---(13)}$$

$$= \frac{N_i^2}{n_i^2}\,\sigma_i^2 \quad \text{---------} \quad (14)$$

The last form arises by neglecting the $=1$ compared with $N_i$, an approximation that will be assumed in what follows and for which no apology need be made when $N_i$ is large.

Interpretation of the result. The equations just written are extremely important; they express the optimum allocation of the sample. For turned around, Eq. 14 can written in the form

$$\frac{n_i}{N_i} = \frac{\sigma_i}{\lambda \sqrt{c_i}} \quad \text{-----------} \quad (15)$$

This equation tells us that if Var x is to be a minimum, the sampling-size in any class should not only be in proportion to the size of that class,

but also in should not only

but also in proportion to the standard deviation of the inventories in that class. Moreover, if the unit cost $c_i$ of collecting information within a particular class is relatively very great, the sample in that class should be cut down because of $c_i$ in the denominator. However, in practice, one must not forget that differential sampling ratios may introduce complications into the field and office procedures. There is simplicity in strict proportionality, wherein $n_i/N_i$ is a constant for all classes, and the statistician should therefore not hastily prescribe differential sampling ratios in accordance with Eq. 15 without some preliminary calculations on gains and costs.

It is usually possible to determine in advance what the costs and sampling errors will be under various plans, and to make a rational choice of sampling ratios By evaluating $\lambda$ it is possible to assign an absolute numerical magnitude to $n_i$. To do this, go back to Eq. 14 and multiply each side by $n_i$, then add for all classes, using the summation sign $\Sigma$ to indicate the operation of summing. The result is

$$\lambda^2 \Sigma\, n_i\, c_i = \frac{\Sigma N_i^2}{n_i}\sigma^2 \qquad \text{But by Eqs. 7 and 15 this gives}$$

$$\lambda^2 C = \lambda \Sigma\, N_i\, \sigma_i \sqrt{c_i}$$

whence

$$\lambda = \frac{\Sigma N_i\, \sigma_i\, \sqrt{c_i}}{C} \quad\text{---------(1b)}$$

Everything on the right is supposed known or approximately determinable beforehand, wherefore a value for $\lambda$ can be calculated and used in Eq. 15 to settle in advance the sample size $n_i$ to be drawn from the various classes. It is possible to make a direct calculation of the minimized value of the Var $X$. This is done by going back to Eq. 4 and rewriting it in the form.

$$\text{Var } X = N_1\left(\frac{N_1}{n_1} - 1\right)\sigma_1^2 + N_2\left(\frac{N_2}{n_2} - 1\right)\sigma_2^2 + \cdots \quad (17)$$

This is obtained from Eq. 4 by dropping the $-1$ in each denominator. As remarked earlier, this approximation is of negligible importance, and the work is greatly simplified.

Going on, the equilibrium value of $N_1/n_1$ is introduced from Eq. 15 with the result that the minimized

$$\text{Var } X = \left(\frac{\lambda \nu c_1}{\sigma_1} - 1\right) N_1\sigma_1^2 + \left(\frac{\lambda \nu c_2}{\sigma_2} - 1\right) N_2\sigma_2^2 + \cdots$$

$$= \lambda \sum N_i \sigma_i \nu c_i - \sum N_i \sigma_i^2 \quad (18)$$

Again, $\lambda$ is to be determined from Eq. 16 then used here. And thus the Var $X$ to be expected from the wisest expenditure of the budget $C$ can be calculated pretty closely beforehand, the only limitation being the prior determinations of $N_1, N_2, \cdots, \sigma_1, \sigma_2, \cdots$. The            on the right arises from the finite multipliers $(N_i - n_i)/(N_i - 1)$ and can be neglected if the sample sizes $n_i/N_i$ are all small.

A simplifying assumption — equal unit costs. The variation in unit costs may be considerable, as when one class of household in the sample is urban and the cost of traveling from one household to another is small while another class is rural in a sparsely populated area and the cost of traveling from

one household to another is considerable. Again, one class might be handled by mail at low cost, while another class requires personal interviews at high cost. However, in many enquiries, such as mail reports from business firms, the unit costs of collec -tion and tabulation are nearly equal in all classes. and as cost enters only through the factor $v c_i$ in the denominator of Eq. 15 anyhow, it is then convenient and satisfactory under such circumstances to put

$$c_1 = c_3 = \cdots \cdots = c \qquad (19)$$

and

$$\qquad (20)$$
$$\lambda \, v_c = K$$

whereupon

$$\frac{n_1}{N_1} = \frac{\sigma_1}{K} \qquad (21)$$

so that

$$\qquad (22)$$
$$\frac{n_1}{N_2} \, ? \, \frac{n_2}{N_3} = \sigma_1 \, ? \, \sigma_2$$

Thus, if $\sigma_1 = 2\sigma_3$, then the sampling interval $N_1/n_1$ in class 1 is only half that of $N_2/n_2$ in Class 2.

- - - - - - - - - - - - - - - - - - - - - -

This interpretation is introduced here for the very practical reason that a systematic selection within classes is permissible and even advantageous in pl -ace of a random selection in many types of sam- pling problems.

If the unit coasts are so nearly equal that Eq.

19 can be assumed, then by solving Eq. 21 for $n_i$ and summing over all class it follows that

$$K = \frac{\sum N_i \sigma_i}{n} \qquad (23)$$

whence upon returning to Eq. 21, we see that it can be written

$$\frac{n_i}{N_i} = \frac{n \sigma_i}{\sum N_j \sigma_j} \qquad (24)$$

wherein $n$, the total sample, must satisfy the relati-ons

$$n = n_1 + n_2 + \cdots\cdots = \frac{C}{c} \qquad (25)$$

Under the assumption of equal unit costs (Eq. 19). Eq. 18 for the minimised $\text{Var } X$ reduces to

$$\text{Var } X = \left( \frac{K}{\sigma_1} - 1 \right) N_1 \sigma_1^2 + \left( \frac{K}{\sigma_2} - 1 \right) N_2 \sigma_2^2$$

$$+ \cdots\cdots\cdots$$

$$= K \sum N_i \sigma_i - \sum N_i \sigma_i^2$$

$$= \frac{(\sum N_i \sigma_i)^2}{n} - \sum N_i \sigma_i^2 \qquad (26)$$

Along with Eqs. 15 or 21, this equation is one of the most important ones in the chapter. It is very handy for calculation, for is evaluates at once the minimised $\text{Var } X$ without the intermediate step of first finding $n_1, n_2, \cdots\cdots$ and substitut-ing them into Eq. 4. As a matter of fact, if it is known that the sampling fraction $n_i/N_i$ are all going to be small, the second term in the last

equation may be neglected, at least for a ready evaluation of the expected precision, for which the from om

$$C.V. \chi = \frac{1}{\sqrt{n}} \frac{\Sigma N_1 \sigma_1}{\xi} \qquad (27)$$

may be preferred. This form has the drawback of requiring an advance estimate of the total inventory $\xi$, which is just the characteristic that it is desired to measure, but this is usually not serious because in practice c.v.x is not needed with extreme accuracy.

The right—hand fraction in the east equation can be written

$$\frac{\Sigma N_1 \sigma_i}{\xi} = \frac{N_1 \sigma_1 + N_2 \sigma_2 + \cdots}{N_1 \mu_1 + N_2 \mu_2 + \cdots} \qquad (28)$$

which can be regarded as a generalized coef
-ficient of variation of the stratified universe.
Eq. 27 then appears as

$$C.V.x = \frac{C.V. universe}{V n}$$

exactly as was encountered in an earlier cha
-pter (p.  ) in consideration of sampling from
an unstratified universe, the finite multiplier
being placed equal to unity.

      If there is only one class in the
universe, the right-hand side of Eq. 28
reduces to $\sigma : \mu$ (dropping the subscript),
and this is the ordinary definition of the
coefficient of variation of a distribution.

      Historical note the advantages of
allocating the sample according to Eq. 21
(or its equivalent, Eq. 27) were first
noted by Neyman. The reference is to p. 5
80 of his article entitled. on the two
different aspects of the representative m-
ethod, journal of the Royal Statistical
Society, Vol. XCV11, 1934; PP. 558-606.

      Remark 1. In sample-design on is
often confronted with the decisilon between

two plans—
    Plan A: Sample-sizes proportional to $N_i$;
    Plan B: Sample-sizes proportional to $N_i\sigma_i$;
                    (Eq. 21 or 24).

One should not too hastily specify Plan B. It
will always show an apparent gain (_vide
infra_) but it is essential to bear in mind
that the variable sampling ratios of Plan B
also require variable weighting factors, and
that the increased costs and complications
in the field-work and tabulation may wipe
out the apparent savings. The expected sav-
ings are not difficult to compute, and do
not require accurate advance estimates of
the standard deviations $\sigma_1$, $\sigma_2$,----within
classes. Preliminary calculations should
always be carried out before making a
decision. Simple examples are given
further on.
      It is often possible to classify the
elements of the universe into strata in
such manner that the standard
deviations $\sigma_i$ are almost sure to be
approximately equal.
Plans A and B are then identical, and
maximum efficiency and simplicity are
attained. This is often accomplished

well enough by grouping the elements on equal class—intervals as of a provious consus.

Remark 2. Every survey is multipurpose, and sometimes two or more characteristics are competitive—that is, the sample design that is adequate for attaining the desired reliability in one important characteristic is not adequate for attaining the required reliability in another important characteristic. Under such conditions the best procedure seems to be strata of approximately equal range (as of the last census), when ordinarily the standard deviations of the competing characteristics will be approximately equal. Proportionate sampl-ing (Plan A) then attains maximum simplicity and approximately maximum efficiency for all the competing characteristics.*

A further simplifying assumption, often usful. In a wide variety of sampling problems it can be assumed that when the elements (dealers, areas) of the universe are grouped according to same measure of size (inventory, sales, number of employees, number of dwelling units) as of a past consus, the standard deviations in the various classes at a later

date will be nearly proportional to their sizes. In symbols,

$$\frac{\sigma_i}{\mu_i} = h, \text{ a constant } \left[\begin{array}{c}\text{First exploited by}\\\text{Hansen and Hurwitz}\end{array}\right] \quad (29)$$

* Jessen, Blythe, Kempthorne, and Deming, "On a population sample for Greece," Journal of the American Statistical Association, vol 42, 1947. This assumption is often admittedly rough, yet extremely useful. Substitution into Eq. 26 gives

$$Var\, x = h^2 \left(\frac{\xi^2}{n} - \sum N_i \mu_i^2\right) \quad (30)$$

in which approximations to $h$ and $\xi$ give a ready value of $Var\, x$.

If the finite multipliers in Eq. 4 are all nearly unity, the second term in parenthesis will be small compared with the first and will serve its purpose even though approximated with still less relative accuracy than $\xi$. In fact, if the second term is neglected altogether there is left the very simple relation

$$C.V.\, x = \frac{h}{\sqrt{n}} \quad (31)$$

wherein $n$, the total sample, comes from Eq. 25.

By introducing Eq. 29 when applicable, the sampling fractions $n_i/N_i$ may be found quickly by noting that

$$\frac{n_i}{N_i} = \frac{n \sigma_i}{\Sigma N_i \sigma_i} \qquad [Eq. 24]$$

$$= \frac{n}{\xi} \mu_i$$

$n$ has disappeared and it is only necessary to assume a mean $\mu_i$ for each class in order to achieve the allocation demanded by Eq. 21. These means ($\mu_i$), should preferably be advance estimates referring to today, rather than to the last census; then today's estimate of $\xi$ is simply $\Sigma N_i \mu_i$ (see the numerical illustrations further on), and $\sigma_i/\mu_i$ in Eq. 29 represents today's coefficient of variation estimated in advance for Class 1.

Actually, in Eqs. 30 and 31 it is only necessary that $n$ be an average value of $\sigma_i : \mu_i$. In using Eq. 32 it is possible to make allowances for departures from the average (see Step iv in the next section).

The assumption contained in Eq. 29 is the constancy of the coefficient of variation of the distribution of inventories, sales, or

employed, from one class to another. It
is a reasonable one in saying that the
variability amongst big dealers is bigger
than the absolute variability amongest
small once. The assumption of strict
proportionality, however, can not be expected
to hold very closely, but fortunately the
usefulness of the assumption is not seriou
-sly impaired if some classes depart only
moderately from strict proportionality,
provided we have some knowledge of
there departures and are able to make
corrections (again see step iv in the next
section). For
example, it is usual to find considerably
greater values of $\sigma_j : \mu_j$ in the smaller
classes, a good example being the blocks
occupied by none or only 1 or 2 families
at the time of a census: a few year
later there may be many new houses
on these blocks where there were only
vacant lote at the time of the census. It
might be well to double the sample of
this class of block, over the sample-size
calculated from Eq. 32. If there is any
danger of a huge housing project being
cuilt meanwhile, any such areas should

be thrown out for a separate sampling —procedure (cf. the book by Hansen, Hurwitz, and Deming mentioned in footnote ).

<u>Steps in the use of the above results in sample-design</u>. The assumption will be made here that the unit costs are the same in all classes, as specified by Eq. 19. Then if the sample is to be allocated according to Eq. 21. the procedure is as follows:

<u>Step 1.</u> Prescribe an allowable coast $c$ and calculate the total allowable sample size $n = C/c$. $c$ is assumed known.

<u>Step ii.</u> Arrive at suitable advance approximations for $\xi$, $h$, and $\mu_i$ if Eq. 29 is useable; otherwise arrive at suitable approximations for $N_i$ and $\sigma_i$. This step requires considerable knowledge of the universe, as might have been acquired in provious surveys or perhaps in a pilot study, as is often a wiss plan.

<u>Step iii.</u> Calculate the minimized expected Var $X$ or C.V. $X$ from Eqs. 26, 27, 30, or 31 (whichever is applicable).

<u>Step iV</u> If the expected precision, just calculated, appears to be good enough, proceed to calculate the sampling fractions $n_i/N_i$ by Eqs. 15, 21, or 32, the latter if applicable, in which case reise or lower sample-sizes by appropriate amounts in those classes wherein the value $h$ assumed for $\sigma_i : \mu_i$ is known to be too low or too high.

If, on the other hand, the expected Var $x$ as computed in Step iii is not considered to be small enough, then it will be necessary to increase the allotted budget from $c$ to some new value $c'$, and recompute. If additional funds are not to be had, it will be necessary either to accept a lower standard of precision (greater Var $x$) or to abandon the survey as not being worth while for the allowable funds.

If the assumption of equal unit is not justifiable, these stens require some obvious modifications.

Solution of the second problem: Var x prescri-bed. Here, C is to be made a minimum, the $n_i$ to be of such magnitude that Var x has a prescribed value (such as C.V. X=2 percent, 5 percent, 25 percent, etc.)

The mathmatical solution is indentical with the first problem down to the point where λ is determined. It is now to be found from Eq. 18. which gives

$$\lambda = \frac{Var\ x + \Sigma N_i \sigma_i^2}{\Sigma N_i \sigma_i V_{c_i}} \qquad (33)$$

Everything on the right is supposed known so λ can be calculated ≢ Then from Eq. 16 the cost is seen to be.

$$C = \frac{1}{\lambda} \Sigma N_i \sigma_i V_{c_i} \qquad (34)$$

The required sample sizes $n_i$ would be compu-ted from Eq. 15 at the same time. The Total sample size will of course be

$$n = \Sigma n_j \quad [\text{As in Eq. 25}]$$

Simplifications may be introduced under the same conditions as those mentioned in earlier paragraphs. In the frist place, if the finite multipliers in Eq. 4 can all be

replace by unity, the second term in the numerator of Eq 33 can be thrown away in the calculation of $\lambda$ Second, if the unit cost are the same in all classes, so that Eq 19 can be used, then in place of Eqs 33 and 34 it would be simpler to write.

$$K = \lambda' c = \frac{Var\, X + \Sigma N_i \sigma_i^2}{\Sigma N_i \sigma_i} \qquad (35)$$

(The second term, arising from the finite multipliers, can sometimes be neglected, as mentioned)

$$and \quad C = \frac{C}{K} \Sigma N_i \sigma_i \qquad (36)$$

The sample sizes ($n_i$) would be computed from Eq 21 at the same time. The total sample will be

$$n = \Sigma n_i = \frac{c}{c} = \frac{\Sigma N_i \sigma_i}{K} \qquad (37)$$

If in addition, the constancy of the ratio $\sigma_i / \mu$, may be usefully assumed, as expressed by Eq. 29, then Eq. 35 gives

$$K = \lambda' \sigma \frac{Var\, X + h^2 \Sigma N_i \mu_i^2}{h \mathcal{E}} \qquad (38)$$

(The second term, arising from the finite multipliers, can sometimes be neglected, as mentioned)

Whereupon the cost simplifies to

$$c = \frac{C}{K} \sum N_i O_i \quad (Eq\ 36)$$

$$= \frac{ch\mathcal{E}}{K} \quad\quad (39).$$

$$= \frac{ch^2 \mathcal{E}^2}{Var\ x} = \frac{ch^2}{(C.V.x)^2}$$

(Neglecting the second term in Eq 38)

The simple sizes would then be found from Eq 32

When the variance is prescribed, and the unite costs $c_i$ are all assumed to be equal, the steps to be applied in allocating the sample according to Eq 21 can be outlined as follows.

Step I Decide on the prescribed value of $Var\ x$, or perhaps, if easier, think in terms of $C.V.x$ (2 percent, 5 percent, 25 percent, etc.)

Step II Arrive at suitable advance approximations for $\mathcal{E}, h,$ and $\mu$, if Eq 29 is useable; otherwise arrive at suitable approximations for $N_i$ and $O_i$. This step requires considerable knowledge of the universe, as might have been acquired in previous surveys or perhaps in a

... study, as is often a wise plan.

Step III. Calculate $\lambda$ from Eq. 33 and then the minimum cost $c$ from Eq. 34. or preferably from Eq. 36 or 39. if applicable.

Step IV. If the cost is apparently not going to be too heavy, proceed to allocate the sample by Eqs. 15, 21, or 32, the latter if applicable, in which case raise or lower the sample-sizes by appropriate amounts in those classes wherein the value $\eta$ assumed for $\sigma_i : \mu_i$ is known to be too low or too high.

If, on the other hand, the expected minimum cost $c$ as computed in Step III is too high, it will be necessary to accept a lower level of precision than initially prescribed, and to recompute the cost. If this plain is not acceptable, the survey must be abandoned, as being too costly.

If the assumption of equal unit costs is not justifiable these steps require some obvious modifications.

Gains compared for different allocation.

Two simple illustrations will help. Let the problem be the determination of the total inventories of 3000 dealers, which have been divided (for simplicity) into just two classes,

as of some previous census. $q_1$ and $\mu_1$ represent advance estimates, as of today. The data are shown in Table 1. It will be assumed that the

<div align="center">Table 1</div>

| class | $N_1$ | $\mu_1$ | $N_1\mu_1$ | $\sigma_1$ | Sample sizes, $n_1$ Plan A $n_1$ proportional to $N_1$ | Plan B $n_1$ proportional to $N_1\sigma_1$ (Eq. 32) |
|-------|-------|---------|------------|-----------|------|------|
| 1 | 1000 | 100 | 100 000 | 90 | 40 | 60 |
| 2 | 2000 | 50 | 100 000 | 45 | 80 | 40 |
| Total | 3 000 | xxx | £=200 000 | xxx | 120 | 120 |

unite ~~costs~~ costs are the same in both classes, wherefore the total costs of both plans are equal because the total sizes are equal. (120). Eq. 4 will give the variance of the total estimated inventory, but for simplicity the finite multipliers will be neglected. Then

$$Var\, x = N_1^2\, \frac{\sigma_1^2}{n_1} + N_2^2\, \frac{\sigma_2^2}{n_2} \quad (\text{Approximation to Eq. 4})$$

Under Plan A this gives

$$Var\, x = 1000^2 \times \frac{90^2}{40} + 2000^2 \times \frac{45^2}{80}$$
$$= \frac{1000^2 \times 45^2}{40}(4+2) = \frac{1000^2 \times 45^2}{120} \times 18$$

Under Plan B the same formula gives

$$\text{Var } x = 1000^2 \times \frac{90^2}{60} + 2000^2 \times \frac{45^2}{60}$$

$$= \frac{1000^2 \times 45^2}{10}(4+4) = \frac{1000^2 \times 45^2}{120} \times 16$$

The loss in using Plan A is thus $(18-16)/16$ or one-eighth. That is, Plan B will produce an eighth less variance at the same cost, or will cost an eighth less for the same variance. Under other conditions (different $N_i$ and $\sigma_i$) the gains arising from Plan B may be much greater, but here the gain is not much more than enough to offset the additional costs and complications in the field and tabulations arising from the introduction of different sampling ratios in the two classes.

This illustration is introduced to show the danger of jumping too rapidly into one plan or the other without some preliminary computations. Even rough guesses at the standard deviations $\sigma_1$ and $\sigma_2$ may be sufficient to indicate which type allocation will best serve the purpose.

The student should satisfy himself that Eq 26 with the second terms neglected gives the same result for Plan B. Thus,

(3/3)

$$Var\ x = \frac{(\sum N_i \sigma_i)^2}{n}$$

$$\{Eq. 26, second\ term\ neglected\}$$
(41)

$$= \frac{(1000 \times 90 + 2000 \times 45)^2}{120}$$

$$= \frac{1000^2 \times 45^2}{120} \times 16 \quad as\ before$$

the
Comparison with the result for $\overline{Var\ x}$ using
no stratification at all is interesting. The
total variance of the universe is

$$\sigma^2 = P\sigma_1^2 + Q\sigma_2^2 + PQ(\mu_2 - \mu_1)^2\ [P. \quad ](42)$$

$$= \frac{1}{3} \times 90^2 + \frac{2}{3} \times 45^2 + \frac{1}{3}\frac{2}{3}(100-50)^2$$

$$= \frac{41,450}{9}$$

P and Q have denote the proportions in the two
class; $P = N_1/N$ and $Q = N_2/N$ where $N = N_2 + N_2$.
An unstratified random sample of 120
will give

$$Var\ x = N^2\frac{\sigma^2}{n}\ [The\ finite\ multiplier\ neglected$$

$$= 3000^2 \times \frac{41,450}{120 \times 9} = \frac{1000^2 \times 45^2}{120} \times 20.5\ (43)$$

It is now possible to compare the variance arising from all three sampling procedures.

Table 2 illustrates a further important point, viz.,

<div align="center">

### Table 2

</div>

| Procedure | $Var\, x$ |
|---|---|
| Unstratified random sampling | $(1000^2 \times 45^2 / 120)\, 20.5$ by Eq 43 |
| Plan A: sample-sizes proportional to $N_i$ | $(1000^2 \times 45^2 / 120)\, 18$ by Eq. 40 |
| Plan B: sample-sizes proportional to $N_i \sigma_i$ | "     $16$ by Eq. 40 or 41. |

That under certain conditions the gains arising from stratified sampling, even when the sample-sizes are proportional to $N_i$ or $N_i \sigma_i$, will be disappointing.

The conditions under which stratification pays off are difficult to describe in general terms, but it is probably safe to say that no striking gains will be made unless the stratification produces widely different means or variance

In the latter case the advantage of using sample-sizes proportional to $N_i \sigma_i$ may be very much worth while. Another illustration (Table 3), being but a simple alteration of the

foregoing one, may make this point clear. The only essential difference between this and the forgoing illustration is that the ratio $\sigma_1 : \sigma_2$ is now 4 instead of 2, yet the results, shown in Table 4, are obviously con͜siderably different from those in Table 2.

In fact, Plan B in Table 4 is seen to require only two-thirds the sample-size of Plan A, which in turn requires only three-quarters the sample size in.

### Table 3

| Class | $N_i$ | $\mu_i$ | $N_i\mu_i$ | $\sigma_i$ | Sample sizes, $n_i$ Plan A $n_i$ Proportional to $N_i$ | Plan B $n_i$ Proportional to $N_i\sigma_i$ (Fig. 3?) |
|---|---|---|---|---|---|---|
| 1 | 1000 | 200 | 200 000 | 2000 | 80 | 160 |
| 2 | 2000 | 50 | 100 000 | 50 | 160 | 80 |
| Total | 3000 | XXX | $\Sigma$ = 300 000 | XXX | 240 | 240 |

### Table 4

| Procedure | Var $\bar{X}$ |
|---|---|
| Unstratified random sample | $(1000^2 \times 50^2/80)$ 24 by Eq. 42 |
| Plan A: sample-sizes proportional to $N_i$ | " 18 by Eq. 40 |
| Plan B: sample-sizes proportional to $N_i\sigma_i$ | " 12 by Eq. 40 or 41. |

unstratified random sampling. In this example, samples proportional to $N_i \sigma_i$ would most certainly show a substantial saving over and above the extra costs of differential sampling and weighting. (The fact that the total sample in this second example is 240 in ~~240 in~~ place of 120 in the first one is of no consequence whatever.)

The calculations are left as an exercise.

As a further exercise the student might calculate $\text{Var } \bar{x}$ using still other sample allocations, such as $n_1 = 200$, $n_2 = 40$; also the reverse.

The student should also calculate $C.V. \bar{x}$ for Plan B in both illustrations by using Eq 31, then verifying the results by calculating $\sqrt{\text{Var } \bar{x} / \bar{x}^2}$, taking $\text{Var } \bar{x}$ and $\bar{x}$ from the table.

Note in regard to advance estimates of the variance. The variance $\sigma_1^2, \sigma_2^2, \cdots$ within the several classes at the time the sample is taken are assumed to be known in advance, approximately at least, as has been explained.

Oftentimes considerable ingenuity is required to evaluate these variances without actually carrying out a full-scale survey. One should make use of census information and other previous studies, consulting with experts in the subject-matter, and

sometimes a pilot study for estimating some of the variance as well as for testing the questionnaire and instructions. These always seems to be a way out if an honest effort is made, particularly in view of certain favorable features of the problem.

For instance, it is important to note that successful allocation does not require accurate knowledge of the variance involved. This is indeed fortunate.

It is so because $Var\ x$ does not increase rapidly from its minimum value as $n_1, n_2, \cdots$ are moved only small distance $\delta n_1, \delta n_2 \cdots$ away from their equilibrium values: this can be seen from Eq. 8 in which the increase in $Var\ X$

arising from the _first_ powers of $\delta n_1, \delta n_2, - - - - - -$
has been placed equal to zero (the condition for
a proper minimum), wherefore Var X is affected
only through the squares and higher powers
of $\delta n_1, \delta n_2, - - - -$ . Now if $\delta n_1$ is small, its
square is still smaller, and the effect on Var X
may be negligible. Hence it is not necessary to
come out with exactly the equilibrium values of
$n_1, n_2, - - - - 3$ approximations will still serve
the purpose tolerably well, which means
that only reasonably good values of $\sigma_1, \sigma_2,$
$- - - - -$ are required.

This point is brought out in the
first numerical illustration, wherein two
different allocations of the sample did not
produce greatly different values of Var X.
In the second illustration the two alloca
-tions of the sample were sufficiently
different to produce distinctly different
values of Var X.

Another noteworthy feature is that optimum
allocation is obtained if only the advance
estimates of $\sigma_1, \sigma_2, - - - -$ are in the _right_
proportions. If all are estimated too high by
(e.g.) 20 percent, as when h in Eq. 29 is
set too high, the sampling fractions $n_i / N_i$ are

unaffected, as is easy to see form Eqs. 24 or 32. Damage is done, however, to the advance estimate of the Var x or the estimated cost. For instance, if the variances were all initially over-estimated by 20 percent, then for a prescribed cost C (Problem I) the advance calculation of the C.V. x will also be too high by about 10 percent, as will be discovered when the variances are evaluated from the returns (cf. the section; "A word on the final determination of the precision reached"). In other words, the error band will turn out to be 10 percent smaller than expected. The same thing will happen when a particular C.V. x has been prescribed (Problem II); the cost will then be about 20 percent greater than was necessary to reach the prescribed precision. The additional expenditure goes for unneeded precision, purchased interestingly enough at the cheapest possible rate (optimum allocation).

Exact formulation of the gains arising from stratification. In the first place, this title implicitly bears with it the thought of allocating the sample in one way or another to the various classes. It is not enough to specify the method of stratification; one must also specify the manner of drawing the sample before he can calculate the variance of the result.

Here we shall compare three plans of allocation, the same three in fact that were just illustrated with numerical examples. Now, however, we shall seek an exact formulation of the comparisons. The variances associated with the three different plans are those shown below.

Plan A: in which $n_i$ is proportional to $N_i$. Here

$$Var \; x = \sum N_i \, (b-1) \, \sigma_i^2$$
$$\left[ \text{From Eq. 17} \right] \quad (44)$$

wherein for all classes

$$\frac{N_i}{n_i} = b \; \left[ \begin{array}{l} \text{the sampling interval} \\ \text{for all classes} \end{array} \right] \; (45)$$

Plan B: in which $n_i$ is proportional to $N_i \sigma_i$. Here

$$Var \; x = \frac{(\sum N_i \sigma_i)^2}{n} - \sum N_i \, \sigma_i^2 \quad [Eq. 26]$$

Plan C: in which there is no stratification at all. Here

$$Var \; x = N^2 \frac{N-n}{N-1} \frac{\sigma^2}{n} = N(b-1)\sigma^2$$
$$\left[ \begin{array}{l} N-1 \text{ replaced} \\ \text{by } N \end{array} \right] \quad (46)$$

wherein

$$\bar{b} = N/n$$

as in Eq. 45

Let these three variances be designated by A, B, C. Then

$$
\begin{aligned}
C - A &= (\bar{b}-1)\left[N\sigma^2 - \sum N_i \sigma_i^2\right] \\
&= (\bar{b}-1)\sum N_i (\sigma^2 - \sigma_i^2) \\
&= N(\bar{b}-1)\sum P_i (\sigma^2 - \sigma_i^2) \\
&= N(\bar{b}-1)(\sigma^2 - \sum P_i \sigma_i^2) \\
&= N(\bar{b}-1)\sum P_i (\mu_i - \mu)^2 \qquad (47)
\end{aligned}
$$

Herein

$$P_i = \frac{N_i}{N}, \text{ the proportion of the universe } (48)$$
$$\text{in Class } i$$

and the line follows from the preceding one from the fact that

$$\sigma^2 = \sum P_i \sigma_i^2 + \sum P_i (\mu_i - \mu)^2 \qquad (49)$$

The relative gain of Plan A over Plan C will be

$$\frac{C-A}{C} = \frac{\sum P_i (\mu_i - \mu)^2}{\sigma^2} \qquad (50)$$

and this relative gain is the same whether it is the _total_ inventory, total sales, total unemployed this is being estimated, or the _average_ inventory or _sales_ per dealer, or the average unemployed per

family or per area. The result shows that stratification · with size of sample proportional to $N_i$ (Plan A) removes the variance between classes.

obviously, this will be appreciably better than straight random unstratified sampling only if one or more of the large strata have widely different means, so that at least one term $N_i (\mu_i - \mu)^2$ will contribute sufficiently to the difference A – c. Thus Plans A and C are compared.

Now compare Planes A and B.

$$A - B = \frac{N}{n} \Sigma N_i \sigma_i^2 - \frac{(\Sigma N_i \sigma_i)^2}{n}$$

$$= \frac{1}{n} \left[ \Sigma N N_i \sigma_i^2 - (\Sigma N_i \sigma_i)^2 \right]$$

$$= \frac{1}{n} \left[ \Sigma N N_i \sigma_i^2 - \Sigma N_i^2 \sigma_i^2 - \Sigma \Sigma N_j N_j \sigma_i \sigma_i \right] \quad i \neq j.$$

$$= \frac{1}{n} \left[ \Sigma N_i \sigma_i^2 (N - N_i) - \Sigma \Sigma N_i N_j \sigma_i \sigma_j \right].$$

$$= \frac{1}{n} \left[ \Sigma N_i^2 \sigma_i^2 + \Sigma \Sigma N_i N_j \sigma_i^2 - \Sigma N_i^2 \sigma_i^2 - \Sigma \Sigma N_i N_j \sigma_i \sigma_j \right]$$

$$= \frac{1}{n} \sum \sum N_i N_j (\sigma_i - \sigma_j)^2 \qquad i < j$$

$$= \frac{1}{2} N \sum \sum P_i P_j (\sigma_i - \sigma_j)^2 \qquad (51)$$

Again, it is probably best to compute $\frac{A - B}{B}$ for the relative gain of Plan B over A. The gain will be worth while only if two or more of the $\sigma_i$ are greatly different, and only then provided neither class is very small. For although $(\sigma_i - \sigma_j)^2$ be large for one pair of classes, the factor $P_i P_j$ may reduce the term to some negligible amount if either Class $i$ or Class $j$ is relatively very small, and the extra cost of using sample sizes proportional to $N_i \sigma_i$ rather than the uniform density of Plan A might not be worth while.

In particular, if two classes have equal variance they contribute nothing at all to the advantage of Plan B over Plan A. To go further, if all classes have the same variance, Plans A and B are identical.

The student should return to the numerical illustration of the last section and test out Eqs. 50 and 51 for the differences between $C - A$ and $A - B$.

A word on the final determination of the precision reached. At the risk of repetition, a word is offered here on the importance of evaluating the precision actually reached, for at least some of the characteristics computed from the sample.* What

(324)

has gone before in this chapter is advice on sample
-design, which in the modern sense implies advance
calculation of the expected variance for some critical
characteristic to be obtained from the sample (total
inventory, sales, or employment). This advance
calculation can only be made on the basis of incompl.
knowledge — the best obtainable at the time, but
incomplete nevertheless. For example, one never
knows what snags the field-work will run into.
but these sangs affect the variances $\sigma_1^2$, $\sigma_2^2$, -, ---.
A subsample take from the returns will give estimates
of these variances as they actually occurred in the
sampling. Substitution into Eq. 4 then gives a close
estimate of Var x.

      It is usually not permissible to use
Eq. 26 for evaluating Var x after the survey is
completed because oftentimes the actual field
procedure and sample-sizes differ consider
-ably from those intended, owing to
special and unforeseen circumstances
encountered, as

Examples: 1. Deming and Simmons, "On the design of a
sample for dealer "inventories," Journal of the American
Statistical Assoc., Vol. 41, 1946:
PP. 16-33. ii. Hansen, Hurwitz, and Deming, A chapter in
Population Sampling (Bureau of the Census, 1947). iii. Jessen,
Blythe, and Deming, "On a population sample for Greece",
Journal of the American Statistical Assoc., Vol., 1947: PP. --.

well as sometimes to misunderstandings or ignoring the
instructions. It is therefore better to use Eq. 4 with
the actual samples $n_1, n_2, \ldots$ and the new
estimates of the variances $\sigma_1^2, \sigma_2^2 \ldots$

The proper size of subsample to take for estimating
the variances from the returns has been discussed
elsewhere (pp. ___ - ___).

This last stage -- the final evaluation of Var x -- be it noted, is
independent of the assumptions that went into the design, however
crude, and hence independent of whether the design was really
a good one.

This last stage is the guarantee of precision, and a good sample
deserves it. The cost it trifling.

A further argument lies in the experience to be gained by
computing some actual variances under the conditions met, in order
to have them available for future sample designs.

## Exercises

1. In sampling the dwelling units of a city for estimating the
proportion colored, it is possible to classify the dwelling units into
two groups, in one of which the proportion colored is .1 and in the
other of which it is .9. Prove that sample sizes proportional to
to $N_i$ are identical with sample sizes proportional to $N_i \sigma_i$; in other
words, Plans A and B are identical. (Hint: $\sigma_i^2 = p_i q_i = p_2 q_2 = \sigma_2^2$)

2. If the two groups of dwelling units of the preceding exercise
are of equal size, then stratified sampling with a uniform
sampling ratio will show a gain of 64 percent over no stratification
at all. This is, only 36 percent as many dwelling units would be
required on Plans A or B as on Plan C, to reach a prescribed precision

3. (a) If the proportions colored in the two classes are .25 and .75,
then the gains of Plans A and B drop to 25 percent. (b) If the
proportions are .6 and .4, the gains in efficiency are only 4 percent;
that is, the gains accomplished by stratification are now almost negligible

此処に発表する Dr.W.E.Deming の論文は昨年同博士及米口統計使節団の一人として来朝されたと き統計数理研究所で講演されたものでその節同博士講究録へ発表を承諾されたもので鏑口後手を入れて送付されたものである。　（編輯者 小川所員）