

ANALYTIC STUDIES OF SURVEY DATA

By: H. O. Hartley, Iowa State University

1. *The problem of group comparisons in sample surveys*

One of the main objectives of a sample survey is the computation of estimates of (say) means and totals of a number of characteristics attached to the units of a population. More often than not, however, the data are also used for what is known as an 'analytic study' or a 'critical analysis' of a survey. Such an analysis usually involves the comparison of means and totals of certain subgroups of the population. To fix the idea, a farm-economic survey in the State of Iowa may have been primarily planned to obtain estimates of totals of numerous farm-economic items such as annual hog sales or annual bushels of corn sealed for the total population of farms in Iowa. In a subsequent analytic study one may then be concerned with the comparison of some of these items for certain subgroups of farms such as 'owner operators' and 'tenant operators'. Such subgroups of the total population have been termed 'domains of study' by the U. N. Subcommittee on Sampling and this term is also used by Yates (1949, 1953) who provides certain formulas for the estimation of their means and variances in the more elementary survey designs.

Although such domains are usually fairly well defined, it will often not be known until after sampling which of the domains any particular unit belongs to. Thus the domains with which an analytic study is concerned are normally not represented in the sample in prescribed fixed sample proportions and the number of sampled units in each domain will itself be a random variable. This is perhaps the most characteristic difference between 'domains' in analytic studies and 'treatment groups' of experiments; but there are others, and the main departures from standard analysis of variance conditions may be summarized under three headings:

1. The number of units in the domains (subgroups) are random variables.
2. The population from which the samples are drawn is finite.
3. Sampling is often not simple random but stratified and/or multistage resulting in correlations of the characteristics of units in the same domain as well as of units in different domains.

To illustrate these points and, at the same time, introduce notation required later we give below three examples of simple survey designs.

Example 1. Simple random sample.

Yates (1949, p. 152) gives data* for a simple random sample of 1/20 Hertfordshire farms. The $n = 125$ sampled farms were, after selection, classified into 7 districts and the number of farms and their total acreages (of crops and grass) are shown in Table 1 below.

TABLE 1

Numbers and total acreages for 125 Hertfordshire farms classified in 7 districts after selection

	District Number							Total (mean)
	1	2	3	4	5	6	7	
No. of farms in sample .	15	8	40	24	4	24	10	125
Total acreage	1,935	1,385	4,851	4,034	335	2,027	547	15,114
Mean acreage	129.0	173.1	121.3	168.1	83.8	84.5	54.7	120.9

* Yates (1949) stresses that these data are for illustration only, see his explanations pg. 30-31.

The notation for the entries in Table 1 are shown in Table 1a below.

TABLE 1a

General notation for random sample of n units classified in k domains after sampling

	Domain No.				Total (mean)
	1	2	j	k	
No. of units in sample .	${}_1n$	${}_2n$	${}_jn$	${}_kn$	n
Total of character . . .	${}_1y$	${}_2y$	${}_jy$	${}_ky$	y
Mean of character . . .	${}_1\bar{y}$	${}_2\bar{y}$	${}_j\bar{y} = {}_jy/n_j$	${}_k\bar{y}$	$\bar{y} = y/n$

If we visualize the whole population as subdivided into the domains we reach the notation set out in Table 1b.

TABLE 1b

General notation for population values of domains from which a random sample was drawn

	Domain No.				Total (mean)
	1	2	j	k	
No. of units in pop. domain	${}_1N$	${}_2N$	${}_jN$	${}_kN$	N
Total of characteristic . .	${}_1Y$	${}_2Y$	${}_jY$	${}_kY$	Y
Mean of characteristic . .	${}_1\bar{Y}$	${}_2\bar{Y}$	${}_j\bar{Y}$	${}_k\bar{Y}$	$\bar{Y} = Y/N$

The main purpose of an analytic study would now consist in estimating the domain means and totals ${}_j\bar{Y}$ and ${}_jY$ from the sample and to provide errors for such estimates. In the present simplest case of a random sample it is easy to guess (as will, in fact, be established later) that the domain population means ${}_j\bar{Y}$ may be estimated by the corresponding sample means ${}_j\bar{y}$. However, the computation of the errors of means of the 'single classification' given in Tables 1 and 1a by standard analysis of variance technique would take no account of either the sampling procedure by which the data were collected or of the actual population for which estimates are required and would, in fact, introduce the assumption of an artificial model not necessarily relevant to the data. That faulty inferences* can be drawn from the application of standard analysis of variance procedures is obvious if we consider the special case when sampling is 100%, i.e. $n = N$ and when the ${}_j\bar{Y} \equiv {}_j\bar{y}$ are in fact estimated without error.

Example 2. Simple stratified sampling (Yates, 1949, p. 154).

Only the first two entries in each cell, the number of units ${}_jn_k$ and the totals ${}_jY_k$ are shown in the example Table 2; the means ${}_jy_k = {}_jY_k/n_k$ are not entered.

If we visualize the whole population of units (farms) as likewise classified by strata and domains the population numbers and totals corresponding to those in the sample of Table 2a would be denoted by capital letters, i.e.

$${}_jN_k, {}_jY_k, {}_j\bar{Y}_k = {}_jY_k/{}_jN_k; \quad N_k, Y_k, \bar{Y}_k = Y_k/N_k;$$

$${}_jN, {}_jY, {}_j\bar{Y} = {}_jY/{}_jN, \quad N, Y, \bar{Y} = Y/N;$$

This simple example shows the occurrence of correlation (say) between units in the same domain: Of the ${}_jn$ units giv-

* Note also discussion in Section 8.

TABLE 2

Numbers of farms and total wheat acreages for a sample of $n = 135$ Hertfordshire farms stratified by 'size-group' and classified by District after selection

Size group stratum (acres)	Population size N_h	District Number							Total
		1	2	3	4	5	6	7	
6-20	519	0 0	1 0	0 0	1 0	0 0	1 0	0 0	3 0
21-50	357	1 0	1 10	2 0	1 17	0 0	1 0	0 0	6 27
51-150	519	3 36	3 40	5 40	5 65	2 0	5 19	3 14	26 214
151-300	400	4 63	5 213	10 270	8 305	5 112	6 140	2 60	40 1163
301-500	215	4 320	10 1074	11 659	12 989	3 234	2 0	1 16	43 3292
501-	51	1 114	4 487	2 315	9 1937	0 0	1 72	0 0	17 2925
	Number	13	24	30	36	10	16	6	135
Total	Wheat acreage	533	1824	1284	3313	346	231	90	7621

TABLE 2a

General notation for sample of n units stratified in L strata and classified in k domains after selection

Stratum number	Domain Number				Total (mean)
		1...	j...	k	
$h = 1$	No. of units	1^{n_1}	j^{n_1}	k^{n_1}	n_1
	Total	1^{y_1}	j^{y_1}	k^{y_1}	y_1
	Mean	$\bar{1}^{y_1}$	\bar{j}^{y_1}	\bar{k}^{y_1}	\bar{y}_1
h	No. of units	1^{n_h}	j^{n_h}	k^{n_h}	n_h
	Total	1^{y_h}	j^{y_h}	k^{y_h}	y_h
	Mean	$\bar{1}^{y_h}$	\bar{j}^{y_h}	\bar{k}^{y_h}	\bar{y}_h
L	No. of units	1^{n_L}	j^{n_L}	k^{n_L}	n_L
	Total	1^{y_L}	j^{y_L}	k^{y_L}	y_L
	Mean	$\bar{1}^{y_L}$	\bar{j}^{y_L}	\bar{k}^{y_L}	\bar{y}_L
Total (Mean)	No. of units	1^n	j^n	k^n	n
	Total	1^y	j^y	k^y	y
	Mean	$\bar{1}^y$	\bar{j}^y	\bar{k}^y	\bar{y}

ing rise to the domain mean \bar{j}^y 'a cluster' of j^{n_h} will be found in stratum h and will usually be positively correlated. The correlation is similar to that found in a two-way table of an analysis of variance with unequal cell frequencies but, of course, cannot here be assumed to have been caused by an additive model.

Example 3. Two stage sample, primaries drawn with replacement and probabilities proportional to size (p.p.s.) and secondaries drawn without replacement and with equal probability.

Table 3 below gives data from a 'consumer preference survey' carried out by the Statistical Laboratory of Iowa State College in the City of Des Moines. The survey was arranged as a stratified two-stage design with 50 city blocks as strata from each of which were sampled 2 'segments' (primaries) each

containing an expected number of 5 households (secondaries). After completion of the survey the data in each stratum were classified by segment number $t = 1, t = 2$ and in 6 different Income Groups $j = 1, 2, \dots, 6$. Actually Table 3 gives the totals for 5 strata-groups combined by pooling the answers for all segments 1 and all segments 2 for the 10 strata in each group. The table shows for each stratum group:

- the number of households in each classification cell (m_{jt} , top line)
- the total number of persons jy_t in households in the j^{th} income group and in segments $t = 1$ and segments $t = 2$ respectively
- the average number of persons per household ($\bar{j}y_t$) for each cell of the two-way classification.

Table 3a shows the notation in the general case of a stratum from which n primaries ($t = 1, 2, \dots, n$) were drawn with the sample results subdivided into k domains ($j = 1, 2, \dots, k$). The stratum subscript h has been omitted from all symbols.

TABLE 3

Number of households (a), number of persons (b) and number of persons per household (c) for 467 households sampled from the City of Des Moines and arranged in 5 'strata groups', 6 'income groups' and 2 segments

Stratum group	Segment		Income group						Total
			1 Less than \$25 per week	2 \$25-\$50 per week	3 \$50-\$75 per week	4 \$75-\$100 per week	5 \$100- \$125 per week	6 More than \$125 per week	
1	1	a	—	7	5	13	14	8	47
		b	—	21	14	42	55	25	157
		c	—	3.00	2.80	3.23	3.93	3.12	3.34
	2	a	2	8	6	12	5	13	46
		b	2	23	16	45	21	50	157
		c	1.00	2.88	2.67	3.75	4.20	3.85	3.41
2	1	a	2	9	13	9	10	13	56
		b	2	20	42	26	37	43	170
		c	1.00	2.22	3.23	2.89	3.70	3.31	3.04
	2	a	2	9	9	13	7	10	50
		b	2	14	20	48	30	34	148
		c	1.00	1.56	2.22	3.69	4.28	3.40	2.96
3	1	a	11	8	15	9	—	4	47
		b	21	24	50	31	—	10	136
		c	1.91	3.00	3.33	3.44	—	2.50	2.89
	2	a	1	8	11	7	5	3	35
		b	1	17	36	26	19	17	116
		c	1.00	2.12	3.27	3.71	3.80	5.67	3.31
4	1	a	4	5	15	5	12	42	119
		b	8	12	47	13	3	36	119
		c	2.00	2.40	3.13	2.60	3.00	3.00	2.83
	2	a	3	2	20	9	2	9	44
		b	6	1	50	26	3	28	115
		c	2.00	2.00	2.50	2.89	1.50	3.11	2.61
5	1	a	8	10	10	9	9	12	58
		b	15	26	37	32	27	45	182
		c	1.88	2.60	3.70	3.56	3.00	3.75	3.14
	2	a	—	7	14	8	8	5	42
		b	—	15	50	34	27	9	135
		c	—	2.14	3.57	4.25	3.38	1.80	3.21

If we visualize the whole population of units as likewise classified by primaries and domains the population numbers and totals corresponding to those in the sample of Table 3a would be denoted by capital letters, i.e. by

$$N; jM_t, jY_t, j\bar{Y}_t = jY_t/M_t; M_t, Y_t, \bar{Y}_t = Y_t/M_t;$$

$$jM, jY, j\bar{Y} = jY_j/M; M, Y, \bar{Y} = Y/M;$$

TABLE 3a

General notation for two stage sample of n primaries $t = 1, 2, \dots, n$ containing respectively m_t sampled secondaries which are classified into k domains after sampling

Primary Number	Domain No.				Total
		1 . . .	j . . .	k	
$t = 1$	No. of units	1^{m_1}	j^{m_1}	k^{m_1}	m_1
	Total	$1y_1$	jy_1	ky_1	y_1
	Mean	$\bar{1y}_1$	\bar{jy}_1	\bar{ky}_1	\bar{y}_1
\vdots					
t	No. of units	1^{m_t}	j^{m_t}	k^{m_t}	m_t
	Total	$1y_t$	jy_t	ky_t	y_t
	Mean	$\bar{1y}_t$	\bar{jy}_t	\bar{ky}_t	\bar{y}_t
\vdots					
n	No. of units	1^{m_n}	j^{m_n}	k^{m_n}	m_n
	Total	$1y_n$	jy_n	ky_n	y_n
	Mean	$\bar{1y}_n$	\bar{jy}_n	\bar{ky}_n	\bar{y}_n
Total	No. of units	1^m	j^m	k^m	m
	Total	$1y$	jy	ky	y
	Mean	$\bar{1y}$	\bar{jy}	\bar{ky}	\bar{y}

The primary index $t = 1, 2, \dots, N$ now runs through the complete population of primary units. If the population consists of L strata the stratum index $h = 1, 2, \dots, L$ will precede the primary index t in all symbols.

2. Estimation of domain totals, variance formulas and variance estimation

We begin with the estimation of domain totals leaving that for domain means for the next section. It is clear that any theory of estimating domain totals or means will have to cover, as a special case, the situation when the 'domain' consists of the total population, and this is, of course, the well known theme dealt with in the literature on sample survey methodology. As is well known, this generally accepted theory of estimation is essentially distribution free, does not accept any particular models, is in fact based on the first two moments only and can be roughly described as 'unbiased estimation with optimum variance properties'.

In the following we shall accept this theory of estimation also for the more general problem of estimating domain totals and means.

We recall here certain basic formulas for the estimation of population totals which may be found in most text books on the subject: Appropriate to any particular design there are three basic formulas pertinent to the estimation of the population total. They are in general notation:

(a) The estimate of the population total Y : —

$$\hat{Y} = \hat{Y}(y_i) \quad [1]^*$$

(b) The population variance formula for \hat{Y} : —

$$\text{Var}(\hat{Y}) = V(y_i) \quad [2]^*$$

* The arguments y_i denote the characteristics attached to the individual units; in stratified sampling they would of course show a double subscript y_{hi} , in multistage sampling a triple subscript y_{hki} , etc.

(c) The estimated variance of \hat{Y} : —

$$\text{var}(\hat{Y}) = v(y_i) \quad [3]^*$$

The multi-variable functions $\hat{Y}(y_i)$, $V(y_i)$ and $v(y_i)$ will depend on the particular design which has given rise to the sample of y_i .*

For example in stratified sampling (*Example 2*) we would have

$$\hat{Y} \equiv \hat{Y}(y_{hi}) = \sum_{h=1}^L N_h \bar{y}_h \quad [4]$$

$$\text{Var}(\hat{Y}) \equiv V(y_{hi}) = \sum_{h=1}^L N_h^2 (1 - n_h/N_h) n_h^{-1} S_h^2 \quad [5]$$

$$\text{var}(\hat{Y}) \equiv v(y_{hi}) = \sum_{h=1}^L N_h^2 (1 - n_h/N_h) n_h^{-1} s_h^2 \quad [6]$$

$$\text{where } S_h^2 = (N_h - 1)^{-1} \sum_{i=1}^{N_h} (y_{hi} - \bar{y}_h)^2 \quad [7]$$

$$s_h^2 = (n_h - 1)^{-1} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2 \quad [8]$$

and these formulas for $L = 1$ (one stratum) would yield as a special case the corresponding formulas for simple sampling (*Example 1*). The appropriate formulas for the two stage design of *Example 3* (primaries p.p.s. equal take m of secondaries)

$$\hat{Y} \equiv \hat{Y}(y_{ti}) = M \bar{y} \quad [9]$$

$$\text{Var}(\hat{Y}) \equiv V(y_{ti}) = \frac{M}{n} \sum_{t=1}^N \frac{M_t}{m} (1 - m/M_t) S_t^2 \quad [10]$$

$$+ \frac{M^2}{n} \sum_{t=1}^N \frac{M_t}{M} (\bar{y}_t - \bar{y})^2$$

$$\text{var}(\hat{Y}) \equiv v(y_{ti}) = \frac{M^2}{n(n-1)} \sum_{t=1}^N (\bar{y}_t - \bar{y})^2 \quad [11]$$

$$\text{where } S_t^2 = (N_t - 1)^{-1} \sum_{i=1}^{M_t} (y_{ti} - \bar{y}_t)^2 \quad [12]$$

$$s_t^2 = (n_t - 1)^{-1} \sum_{i=1}^{n_t} (y_{ti} - \bar{y}_t)^2 \quad [13]$$

In general \hat{Y} and $v(\hat{Y})$ will depend on the sampled y_i only, and these are the formulas used in practical evaluation of the estimate and its variance, whilst $V(y_i)$ depends on the entire population of y_i values and is used for the theoretical evaluation of the merits of the estimator. The three formulas are applicable to *any* set of characteristics y_i attached to the population units. This enables us to use these same formulas for the estimation of domain totals and their variances: — To this end we introduce the following characteristics y_i which we shall attach to *all* units in the population.

$$_j y_i = \begin{cases} y_i & \text{if the } i\text{th unit belongs to } j\text{th group (domain)} \\ 0 & \text{otherwise} \end{cases} \quad [14]$$

Now the group total $_j Y$ of our j^{th} group is seen to be the population total of the $_j y_i$ and standard sample survey theory, therefore, provides the following estimators:

(a) The estimate of the domain total $_j Y$: —

$$_j \hat{Y} = \hat{Y}(_j y_i) \quad [15]$$

(b) The population variance formula for $_j \hat{Y}$: —

$$_j V = V(_j y_i) \quad [16]$$

(c) The estimated variance of \hat{Y} :—

$${}_iV = v({}_iY) \quad [17]$$

together with the assurance of unbiasedness

$$E({}_i\hat{Y}) = {}_iY \quad [18]$$

$$E({}_iV) = {}_iV_i \quad [19]$$

The meaning of formulas [15] to [17] is simply that the standard formulas for estimation [4] to [13] be applied to the characteristics [14]. The 'spelling out' of these formulas often results in simplifications: In Example 2 (stratified sampling) we find that formulas [15] to [17] can be written as

(a) The estimate of the domain total:

$${}_i\hat{Y} = \sum_h N_h {}_iY_h/n_h \quad [20]$$

(b) The population variance formula for \hat{Y}

$${}_iV = \sum_h \frac{N_h(N_h - n_h)}{n_h(N_h - 1)} \left\{ \sum_{i=1}^{N_h} {}_iY_h^2 - {}_iY_h^2/N_h \right\} \quad [21]$$

(c) The estimated variance of \hat{Y}

$${}_iV = \sum_h \frac{N_h(N_h - n_h)}{n_h(n_h - 1)} \left\{ \sum_{i=1}^{n_h} {}_iY_h^2 - {}_iY_h^2/n_h \right\} \quad [22]$$

The last formula may be compared with one given in the 2nd edition of Yates' book (Yates 1953, 301 formula [9.3.e]) which corrects an earlier (faulty) formula given in his 1st edition (Yates 1949, 202). To make this comparison we use the sample variance of the units in the h th stratum and in the j th domain, we write

$${}_iS_h^2 = \sum_{i=1}^{n_h} ({}_iY_h - \bar{{}_iY_h})^2 / (n_h - 1)$$

and hence obtain [22] in the form

$${}_iV = \sum_h \frac{N_h(N_h - n_h)}{n_h(n_h - 1)} \left\{ (n_h - 1) {}_iS_h^2 + \left(\frac{1}{n_h} - \frac{1}{N_h} \right) {}_iY_h^2 \right\} \quad [23]$$

which agrees with Yates' formula [9.3.e]. Our formulas [22] and [23] have been proved to be unbiased estimates of the exact variance of \hat{Y} . The second term of [23] which is proportional to the square of the domain total ${}_iY_h^2$, is characteristic of estimators of totals of a population domain whose size is unknown and this feature will, of course, disappear in the estimation of domain means.

In Example 3, (two stage sampling with primaries drawn p.p.s.) if we assume that an equal take of m secondaries are sampled from each primary, formulas [15] and [17] spell out as follows:

(a) The estimate of the domain total:

$${}_i\hat{Y} = M {}_iY/m = M \left(\frac{{}_iY}{m} \right) \bar{{}_iY} \quad [24]$$

(c) The estimated variance of \hat{Y}

$${}_iV = \frac{M^2}{n(n-1)} \sum_{i=1}^n \left(\frac{{}_iY_i}{m_i} - \frac{{}_iY}{m} \right)^2 \quad [25]$$

The interpretation of [24] is simple: The sample domain mean $\bar{{}_iY}$ estimates the corresponding population mean ${}_iY$ and the fraction m/m the corresponding population fraction M/M , the unbiasedness of the product being assured by [18].

3. Estimation of domain means, variance formulas and variance estimation

In the preceding section we gave estimators \hat{Y} of the population total Y of the j th domain. We now turn to the estimation of the domain means \bar{Y} (given by Y/N in the notation of Example 1 and 2 and Table 1 b, and by Y_i/M in the notation of Example 3). Even if the number of units in the j th domain N (or M) were known it would usually be unwise to use \hat{Y}/N (or in two-stage sampling \hat{Y}/M) as an estimator of \bar{Y} , since it is well known from experience with such estimators of population means that their variances are large. In any case N or M will often be unknown (as are, for example, the number of tenant-operators in Iowa) and in these circumstances we are forced to estimate both the unknown numerator Y and unknown denominator N (or M). This direct approach automatically leads to what is known as the "combined ratio estimator" of \bar{Y} . Other ratio estimators, available in special situations, are briefly discussed in section 6. In order to estimate, then, the denominator N (or M) by precisely the same method used for the numerator Y we introduce the "count variates"

$${}_i u_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ unit is in } j^{\text{th}} \text{ domain} \\ 0 & \text{otherwise} \end{cases} \quad [26]$$

and these obtain the estimate of the domain size N (or M) as

$$\hat{U} = \hat{Y}({}_i u_i). \quad [27]$$

Provided that $\hat{U} > 0$, that is provided there is at least one sampled unit in the j th domain, we can then use

(a) the combined ratio estimator of the domain mean \bar{Y}

$$\hat{{}_iY} = \hat{Y} / \hat{U} \quad [28]$$

In general this estimator will be biased and conditions as to when this bias is negligible are known from the literature (see e.g. Cochran, 1953, pp. 114-8). Unbiased ratio estimators are briefly discussed in section 6. Standard approximate formulas (see e.g. Cochran, 1953, pp. 114-120) for the variance and estimated variance of the ratio estimator likewise yield

(b) The approximate variance formula for $\hat{{}_iY}$

$$\text{Var}(\hat{{}_iY}) \doteq {}_iN^{-2} V({}_iY_i - \bar{{}_iY} {}_i u_i) \quad [29]$$

(c) The approximate estimate of the variance of $\hat{{}_iY}$

$$\text{var}(\hat{{}_iY}) = \hat{{}_iU}^{-2} v({}_iY_i - \hat{{}_iY} {}_i u_i) \quad [30]$$

The spelling out of these formulas [28], [29] and [30] will show that the domain mean \bar{Y} is usually estimated by the simple sample domain mean $\bar{{}_iY}$. This will certainly be so in Examples 1 and 3 as indeed with all self-weighting survey designs. Likewise the estimated variance of $\hat{{}_iY}$ will often simplify:

In Example 2 we have the following formulas

$$(a) \quad \hat{{}_iY} = \sum_h \frac{N_h}{n_h} {}_iY_h / \sum_h \frac{N_h}{n_h} {}_i u_h \quad [31]$$

which, in the case of proportional allocation $N_h/n_h = N/n$ will reduce to

$$(a) \quad \hat{{}_iY} = \bar{{}_iY}. \quad [32]$$

The estimated variance in this example from [30] and [23] is given by

$$\text{var}(\hat{{}_iY}) = \hat{{}_iU}^{-2} \sum_{h=1}^L \frac{N_h(N_h - n_h)}{n_h(n_h - 1)} \left\{ {}_i\Sigma_h + {}_iQ_h {}_i u_h (\bar{{}_iY}_h - \bar{{}_iY})^2 \right\} \quad [33]$$

where ${}_i\Sigma_h = ({}_i n_h - 1) {}_iS_h^2$ is the sum of squares of deviations of

the y in the j^{th} domain of the h_{jk} stratum about their mean \bar{y}_{jk} and $i q_k = (n_k - n_{jk})/n_k$ the proportion of units sampled outside the j^{th} domain. Formula [33] agrees with Yates' (1953, p. 301 formula [9.3.c]) except that the latter appears to be restricted to a proportional allocation of the sample to strata. The formula for Example 1 (simple sampling) is obtained as a special case of [33] when $L = 1$ and $\bar{y}_{jk} = \bar{y}$ yielding

$$\text{var } \hat{y} = \left(1 - \frac{n}{N}\right) \left(\frac{n}{n-1}\right) \left(\frac{n-1}{n}\right) s^2/n \quad [34]$$

which is, approximately, equal to the familiar variance of a mean of a sample of size n .

It should be noted that the case of stratified sampling is not derivable from the simple sampling case by summation over strata since this would only yield the first terms inside the $\{\}$ of [33]. The second terms allow for the fact that the strata proportions n_{jk}/N for the j^{th} domain are not known.

A numerical example:

We illustrate the above formulas by evaluating \hat{y} and $\text{var } (\hat{y})$ for Example 2 using the data of Table 2 to evaluate an estimate of the mean wheat acreage for District $j = 1$.

The N_k and n_k are given in the second and last (top entry) columns of Table 2 and the y_k and m_k in the third column. For the individual wheat acreage values y_{ki} of the farms in the first district reference is made to Yates (1949, p. 154). We compute

$${}_1\hat{Y} = 3291, \quad {}_1\hat{U} = 183, \quad \hat{y} = 17.98, \quad \text{var } (\hat{y}) = 61.52.$$

Finally the spelling out of formulas [30] and [28] for Example 3 yields

$$\text{var } (\hat{y}) = \frac{n}{n-1} \sum_{i=1}^n \left(\frac{m_i}{m}\right)^2 (\bar{y}_i - \hat{y})^2 \quad [35]$$

where $\hat{y}_i = \bar{y}_i$.

4. Correlation between the estimates of domain totals and means and variances of differences between means

In the preceding sections we derived formulas for the estimation of domain totals and means and their variances. We now turn to the task of the "comparison" of two estimated domain means, and to this end require estimates of their covariances in order to obtain variance estimates for the difference of two estimated domain means.

Again, we first deal with the estimates of domain-totals and afterwards with the domain means. It is convenient to write all formulas for the comparison of two particular domains which, without loss of generality, may be taken as domains $j = 1$ and $j = 2$.

The following formulas are obvious by arguments similar to those of section 2.

$$\text{Var } ({}_1\hat{Y} - {}_2\hat{Y}) = V ({}_1y_i - {}_2y_i) \quad [36]$$

$$2 \text{ Cov } ({}_1\hat{Y}, {}_2\hat{Y}) = V ({}_1y_i) + V ({}_2y_i) - V ({}_1y_i - {}_2y_i) \quad [37]$$

$$\text{var } ({}_1\hat{y} - {}_2\hat{y}) = v ({}_1y_i - {}_2y_i) \quad [38]$$

$$2 \text{ cov } ({}_1\hat{y}, {}_2\hat{y}) = v ({}_1y_i) + v ({}_2y_i) - v ({}_1y_i - {}_2y_i) \quad [39]$$

Formula [38] gives directly the expression required for estimating the variance of the difference between two estimates ${}_1\hat{y}$ and ${}_2\hat{y}$ of the two domain totals. Formula [39], however, is required for proving that $\text{cov } ({}_1\hat{y}, {}_2\hat{y})$ is given by the bilinear form in the variates ${}_1y_i$ and ${}_2y_i$ which corresponds to the symmetric quadratic form $v(y_i)$. The same remarks apply to the variance formulas [37] and [36].

Turning now to the domain means we require an approximate expression for the covariance between two ratios y/u and x/v : This follows on identical lines as the familiar variance formulas and yields:

$$\text{Cov} \left(\frac{y}{u}, \frac{x}{v} \right) = (E(u) E(v))^{-1} \text{Cov} \left(y - \frac{E(y)}{E(u)} u, x - \frac{E(x)}{E(v)} v \right) \quad [40]$$

The combination of [36] and [39] can be shown to yield

$$\begin{aligned} \text{Var } ({}_1\hat{y} - {}_2\hat{y}) &\equiv \text{Var} \left(\frac{{}_1\hat{Y}}{{}_1\hat{U}} - \frac{{}_2\hat{Y}}{{}_2\hat{U}} \right) \\ &\equiv V \left(\frac{{}_1\hat{Y}}{{}_1\hat{U}} - \frac{{}_2\hat{Y}}{{}_2\hat{U}} \right) \end{aligned} \quad [41]$$

The corresponding approximation for the estimated variance yields

$$\text{var } ({}_1\hat{y} - {}_2\hat{y}) \doteq v \left(\frac{{}_1\hat{Y}}{{}_1\hat{U}} - \frac{{}_2\hat{Y}}{{}_2\hat{U}} \right) \quad [42]$$

where the variables ${}_1y_i$ and ${}_2y_i$ are given by [14] and [26] respectively, the estimates ${}_1\hat{U}$ and ${}_2\hat{U}$ by [27] and [28] and the multi-variable function v is defined in [3]. This function v will of course depend on the particular survey design but there is no difficulty in spelling it out in any particular case. We proceed to do so for our examples 1, 2 and 3 when $v(y_i)$ is given by equations (6) and (11) respectively. For example 2, simple stratified sampling with varying sampling fractions, we obtain the formula

$$\begin{aligned} v ({}_1\hat{y} - {}_2\hat{y}) &= \sum_k a_k \left\{ \frac{{}_1\sum_k}{{}_1\hat{U}^2} + \frac{{}_2\sum_k}{{}_2\hat{U}^2} + 2 n_k {}_1p_k {}_2p_k \frac{({}_1\bar{y}_k - \hat{y})({}_2\bar{y}_k - \hat{y})}{{}_1\hat{U} {}_2\hat{U}} + \right. \\ &\quad \left. + n_k {}_1p_k {}_1q_k ({}_1\bar{y}_k - \hat{y})^2 + n_k {}_2p_k {}_2q_k \frac{({}_2\bar{y}_k - \hat{y})^2}{{}_2\hat{U}^2} \right\} \end{aligned} \quad [43]$$

where

$$\begin{aligned} a_k &= \frac{N_k (N_k - n_k)}{n_k (n_k - 1)}, \quad {}_1p_k = \frac{n_k}{n}, \quad {}_1q_k = 1 - {}_1p_k \\ {}_1\sum_k &= \sum_{i=1}^{n_k} ({}_1y_{ki} - \bar{y}_k)^2 \end{aligned}$$

Formula [43] is in agreement with that for $\text{cov } ({}_1\hat{y}, {}_2\hat{y})$ given by Yates (1953, p. 301, formula 9.3.d) except that the latter appears to be restricted to stratified sampling with proportional allocation. The first two terms in [43] are the within domain within strata components, whilst the last three terms have the form of multinomial variances and covariances and allow for the fact that the strata proportions ${}_1N_k/{}_1N$ and ${}_2N_k/{}_2N$ are unknown. These latter terms disappear in the special case of a single stratum $L = 1$ which yields the answer for example 1 in the form

$$V ({}_1\hat{y} - {}_2\hat{y}) = \left(1 - \frac{n}{N}\right) \left(\frac{n}{n-1}\right) \left\{ \frac{{}_1s^2 - 1}{{}_1n} + \frac{{}_2s^2 - 1}{{}_2n} \right\} \quad [44]$$

This expression is approximately equal to $({}_1s^2/{}_1n + {}_2s^2/{}_2n)$ i.e. the familiar variance formula for the difference of two means based on fixed samples of size ${}_1n$ and ${}_2n$.

Finally, the spelling out of formula [42] for Example 3 (two stage sampling primaries drawn p.p.s. equal take of secondaries) yields the surprisingly simple formula

$$v ({}_1\hat{y} - {}_2\hat{y}) \doteq \frac{n}{n-1} \sum_{i=1}^n \left\{ \frac{{}_1m_i}{{}_1m} ({}_1\bar{y}_i - \hat{y}) - \frac{{}_2m_i}{{}_2m} ({}_2\bar{y}_i - \hat{y}) \right\}^2 \quad [45]$$

In the important case of $n = 2$ primaries this formula simplifies further to

$$v(y - \bar{y}) = 4(w(\bar{y}_1 - \bar{y}_2) - \bar{y}(\bar{y}_1 - \bar{y}_2))^2 \quad [46]$$

where

$$w = m_1 m_2 / (m_1 + m_2)^2$$

Certain special cases of formulas [44] and [45] were recently obtained by an independent derivation by L. Kish and Irene Hess (1955).

5. *The domain total expressed as a proportion of the population total.*

In many analytic studies of sample survey data we are interested in the proportion of a measured characteristic which falls into a specified domain. For example in a consumer study we may be interested in the proportion of say the total milk consumption which is attributable to families of a particular income group. Or, again, in a soil survey we may be interested in the proportion of the total farm land which is of a particular soil type. In the preceding sections we have estimated the total ${}_jY$ of the variate values y_i which fall into the j^{th} domain by the estimator ${}_j\hat{Y}$. It is therefore suggested that the proportion ${}_jY/Y$ be estimated by the ratio

$${}_j\hat{p} = {}_j\hat{Y}/\hat{Y} \quad [47]$$

The properties of this ratio estimator can be evaluated by a method similar to that used in Section 3: We attach two variates to each unit in the population viz the numerator variable given (by [14] i.e.) by

$${}_jy_i = \begin{cases} y_i & \text{if the } i^{\text{th}} \text{ unit belongs to the } j^{\text{th}} \text{ group} \\ 0 & \text{otherwise} \end{cases} \quad [48]$$

and the denominator variable given by y_i . Since ${}_j\hat{Y}$ and \hat{Y} are, respectively, the estimates of the population totals for the variates ${}_jy_i$ and y_i , the estimator ${}_j\hat{p}$ is seen to be the standard (combined) ratio estimator for ${}_jP = {}_jY/Y$ of the respective population totals. Its variance may therefore be obtained by the standard (approximate) formula viz.

$$\text{Var } {}_j\hat{p} = Y^{-2} V({}_jy_i - {}_jP y_i) \quad [49]$$

where the variance function $V(y)$ is defined by [2]. Let us spell out the general formula [49] for the particular case of a simple stratified design and in terms of the notation used in Table 2a and that described immediately following Table 2a. We obtain the formula

$$\text{Var } {}_j\hat{p} = Y^{-2} \sum_h \frac{N_h^2}{n_h(N_h - 1)} \sum_{i=1}^{N_h} (d_{hi} - \bar{D}_h)^2 \quad [50]$$

where $d_{hi} = {}_jy_{hi} - {}_jP y_{hi}$ and \bar{D}_h is the stratum mean of the d_{hi} .

Yates (1953) p. 304 only discusses this case of a stratified sample and offers a formula only for the special case when the domains do not cut across the strata. In the general case of domains cutting across the strata Yates gives hints for variance computation and we have not attempted to identify the result of these instructions with our simple formula [50].

For the purpose of variance comparison, an alternative form of [50] may be more useful. After some algebra we reach the formula

$$\begin{aligned} \text{Var } ({}_j\hat{p}) = Y^{-2} \sum_h \frac{N_h^2}{n_h(N_h - 1)} \{ & {}_jQ^2 {}_j\Sigma_h + {}_jP^2 {}_j\Sigma_h \\ & + N_h^{-1} {}_jN_h {}_jN_h ({}_jQ {}_j\bar{Y}_h + {}_jP {}_j\bar{Y}_h)^2 \} \end{aligned} \quad [51]$$

where

${}_jN_h$ is the number of units in the h^{th} stratum falling within the j^{th} domain

${}_jN_h$ is the number of units in the h^{th} stratum not falling within the j^{th} domain

${}_jY_h$ is the y -total for the h^{th} stratum of units in the j^{th} domain

${}_jY_h$ is the y -total for the h^{th} stratum of units not in the j^{th} domain.

So that

$${}_j\bar{Y}_h = {}_jY_h / {}_jN_h \quad {}_j\bar{Y}_h = {}_jY_h / {}_jN_h$$

are the corresponding means.

Further

$${}_jP = \sum_h {}_jY_h / Y \quad {}_jQ = \sum_h {}_jY_h / Y = 1 - {}_jP$$

$${}_j\Sigma_h = \sum_{i=1}^{N_h} ({}_jy_{hi} - {}_j\bar{Y}_h)^2$$

is the sum of squares of deviations of the y values in the h^{th} stratum and j^{th} domain from their mean and

$${}_j\Sigma_h = \sum_{i=1}^{N_h} ({}_jy_{hi} - {}_j\bar{Y}_h)^2$$

is the corresponding sum of squares for the y -values in the h^{th} stratum and *not* in the j^{th} domain.

To interpret the three terms in [51] we may write the estimator in the form

$${}_j\hat{p} = 1 + ({}_j\hat{Y} / \hat{Y}) \quad [52]$$

where ${}_j\hat{Y} = \hat{Y} - {}_j\hat{Y}$ is an estimate for the y -total *not* in the j^{th} domain. It can be shown that the first two terms of [51] are contributed by variation in y -values only, holding the ratios n_h/n_h constant at their expected values $\frac{{}_jN_h}{N_h} n_h$. The third term allows for the variability in the proportion n_h/n_h and is seen to be a binomial type of variance being approximately equal to

$$\sum_h \left(\frac{{}_jN_h}{N_h} \right) \left(\frac{{}_jN_h}{N_h} \right) \frac{1}{n_h} Y^{-2} N_h^2 ({}_j\bar{Y}_h {}_jP + {}_j\bar{Y}_h {}_jQ)^2. \quad [53]$$

In this equation [53] the term

$$\left(\frac{{}_jN_h}{N_h} \right) \left(\frac{{}_jN_h}{N_h} \right) \frac{1}{n_h}$$

is the binomial variance of n_h/n_h and the term inside the brackets $\{ \}$, is the square of the expected value of the coefficient of n_h/n_h in the expansion of the estimator ${}_j\hat{Y}/\hat{Y}$.

For the computation of an estimate of variance we, again, rely on standard formulas for ratio estimators. The standard estimate (not necessarily unbiased) of the variance given by [51] is computed from

$$\hat{\text{var}} ({}_j\hat{p}) = \hat{Y}^{-2} v({}_jy_i - {}_jP y_i) \quad [54]$$

where the sample variance function is given by [3].

If we spell out this general formula in the particular case of a stratified sample we obtain

$$\text{var } {}_j\hat{p} = \hat{Y}^{-2} \sum_h \frac{N_h^2}{n_h(n_h - 1)} \sum_{i=1}^{N_h} (d_{hi} - \bar{d}_h)^2 \quad [55]$$

where $d_{hi} = {}_jy_{hi} - {}_jP y_{hi}$ and \bar{d}_h is their stratum sample means and \hat{Y} is the estimate of the population total i.e. $\hat{Y} = \sum_h N_h \bar{y}_h$ and ${}_j\hat{p} = {}_j\hat{Y}/\hat{Y}$ the estimate of the proportion. An alternative

formula which does not require the computation of the d_h can be obtained, after some algebra as follows.

$$\text{var } i\hat{p} = \hat{Y}^{-2} \sum_h \frac{N_h^2}{n_h(n_h - 1)} \{ i\hat{q}^2 i\hat{T}_h + i\hat{p}^2 j\hat{T}_h + n_h^{-1} i\hat{n}_h j\hat{n}_h (i\hat{q}_i \bar{y}_h + i\hat{p} j\bar{y}_h)^2 \}. \quad [56]$$

Here

$i\hat{n}_h$ is the number of sampled units in the h^{th} stratum falling within the j^{th} domain.

$j\hat{n}_h$ is the number of sampled units in the h^{th} stratum not falling within the j^{th} domain.

$i\hat{y}_h$ is the y -total for the h^{th} stratum of units in the j^{th} domain.

$j\hat{y}_h$ is the y -total for the h^{th} stratum of units not in the j^{th} domain.

So that

$$\bar{i\hat{y}}_h = i\hat{y}_h / i\hat{n}_h \quad \bar{j\hat{y}}_h = j\hat{y}_h / j\hat{n}_h$$

are the corresponding means. Further

$$i\hat{p} = i\hat{Y} / \hat{Y} \quad , \quad i\hat{q} = 1 - i\hat{p}$$

$$i\hat{T}_h = \sum_{i=1}^{i\hat{n}_h} (i\hat{y}_{hi} - \bar{i\hat{y}}_h)^2$$

is the sum of squares of deviations of the sampled y -values in the h^{th} stratum and j^{th} domain from their mean and

$$j\hat{T}_h = \sum_{i=1}^{j\hat{n}_h} (j\hat{y}_{hi} - \bar{j\hat{y}}_h)^2$$

is the corresponding sum of squares for the y -values in the h^{th} stratum and not in the j^{th} domain.

6. Domain means adjusted for concomitant variables.

We now turn to the analogue to «analysis of variance and covariance» for domain means. First let us briefly recall the essentials of «analysis of covariance» as it is practiced with the adjustment of (say) two group means from experimental designs. It is usually assumed that:

- (a) The true (or population) means for the concomitant variable x are identical for the two groups, so that any differences in observed x -group means are due to sampling. In popular language the failure of this condition to be satisfied is sometimes described by the warning that «Analysis of Covariance should not be misused to correct away real treatment differences in the x -means.»
- (b) The true (or population) regression lines for the two groups are parallel.

Condition (b) is, of course, not essential, but in the situations in which it is satisfied (or approximately satisfied) it is almost universally invoked and simplifies the analysis. If condition (a) is not satisfied a meaningful generalization of analysis of covariance can still be employed if the two population group x -means $i\bar{X}$ are known, and estimates of the y -group means can be obtained as the ordinates of the respective estimated group regression lines evaluated at the respective abscissa values $i\bar{X}$. Now with situations as are usually encountered in «analytic studies» of sample survey data one will hardly ever be able to assume identity (or even approximate identity) of the x -domain means $i\bar{X}$; on the other hand, there are situations when

these means are known and we shall therefore here attempt to develop an analogue to analysis of covariance in this case.

We therefore deal with the following situation: A sample survey provides paired data $y_i x_i$ for a sample of n units sampled from a population of N units. After sampling the n units are classified into k domains $j = 1, 2, \dots, k$ for which the population means $i\bar{X}$ are known. It is now required to estimate the (unknown) domain y -means $i\bar{Y}$ utilizing the $i\bar{X}$.

The question then arises as to which estimator of $i\bar{Y}$ should be used. The regression theory in classical analysis of covariance arises as the maximum likelihood estimation and results from the assumption of a linear model. Although the validity of such a model, even for very large finite populations, will often be in doubt, the use of regression estimation may still result in a gain of precision. Nevertheless we shall here *not* employ regression estimators. The reason for this is *not* that we consider regression theory inappropriate, but that this theory for finite populations requires considerable development before it can be applied in the present situation. On the other hand, ratio estimators are easily adapted to the estimation of domain means but in using these we should stress their well known limitation, namely, that they are likely to be effective only if the y and/or x scales can be so chosen that the population regression will intersect near the origin.

If a combined ratio estimator is used to estimate $i\bar{Y}$ with the help of the known x -domain mean $i\bar{X}$ the theory is almost identical with that developed for the ratio estimation of Section 3, the only essential difference is that the concomitant variable x will now take the place of the «count variable» $i\mu_i$ of [26]. Accordingly we introduce the variate

$$ix_i = \begin{cases} x_i & \text{if } i^{\text{th}} \text{ unit is in } j^{\text{th}} \text{ domain} \\ 0 & \text{otherwise} \end{cases} \quad [57]$$

and obtain the ratio estimator of $i\bar{Y}$ in the form:

- (a) The combined ratio estimator of the domain mean $i\bar{Y}$:

$$\tilde{i\hat{y}} = i\bar{X} (\hat{i\hat{y}} / \hat{i\hat{x}}) \quad [58]$$

where the $i\bar{X}$ are the fixed and known domain means of x , $i\hat{Y}$ is defined by [15] and $i\hat{X}$ by analogy.

Employing standard results from the ratio estimator theory on similar lines as in section 3 we further obtain

- (b) The approximate variance formula of $\tilde{i\hat{y}}$

$$\text{Var}(\tilde{i\hat{y}}) = iN^{-2} V(i\hat{y}_i - i\bar{Y} ix_i / i\bar{X}) \quad [59]$$

- (c) The approximate estimate of the variance of $\tilde{i\hat{y}}$

$$\text{var}(\tilde{i\hat{y}}) = i\hat{U}^{-2} v(i\hat{y}_i - i\hat{Y} ix_i / i\hat{X}) \quad (60)$$

where the multivariable functions V and v are defined by [2] and [3] the variates $i\hat{y}_i$ and ix_i by [14] and [57], $i\hat{Y}$ and $i\hat{U}$ by [15] and [27] and $i\hat{X}$ by analogy.

Formulas [59] and [60] show the familiar result that the variance of $\tilde{i\hat{y}}$ depends on the residuals $i\hat{y}_i - (i\bar{Y} / i\bar{X}) ix_i$ of a «regression» with slope $i\bar{Y} / i\bar{X}$ passing through the origin. For the majority of units in the population $i\hat{y}_i = ix_i = 0$ and no contribution is made to the residuals; for the units in the j^{th} domain the above residuals will be small only if the y, x data of this domain alone satisfy the usual conditions required for the effectiveness of ratio estimators, namely that the y, x correlation should be high and that the y, x regression should intersect near the origin. If the latter condition is not satisfied little is gained by using regression estimators. For if the y, x data in the j^{th} domain have a regression which does not intersect near

the origin then the addition of the large number of units (outside the j^{th} domain) for which $y_i = x_i = 0$ would generate large residuals from the best regression fitted to all N pairs of y_i, x_i .

The spelling out of formula (60) follows on the same lines as in section 3. For example 1 (random sampling) we find in analogy to (34)

$$\text{var}(\tilde{y}) = \left(1 - \frac{n}{N}\right) \frac{n}{j n^2 (n-1)} \sum_{i=1}^n \left(y_i - \frac{y}{j x} x_i\right)^2 \quad [61]$$

In a similar manner the formulas for the variance of a difference of two ratio estimators can be obtained following on the lines of the arguments in section 4. We reach the general formula

$$\text{var}(\tilde{y}_1 - \tilde{y}_2) = v \left(\frac{1}{U} \left(y_1 - \frac{1}{1X} x_1 \right) - \frac{1}{2U} \left(y_2 - \frac{2}{2X} x_2 \right) \right) \quad [62]$$

which, in the case of example 1 (simple sampling) spells out as

$$\text{var}(\tilde{y}_1 - \tilde{y}_2) = \left(1 - \frac{n}{N}\right) \left(\frac{n}{n-1} \right) \left\{ \frac{1}{1n^2} + \frac{2}{2n^2} \right\} \quad [63]$$

where

$$j \Sigma = \sum_{i=1}^n \left(y_i - \frac{y}{j x} x_i \right)^2$$

7. Unbiased ratio estimators for domain means and exact variance formulas.

In sections 3-6 we have freely used analogies to the well known combined ratio estimators and their approximate variance formulas. For a discussion of the magnitude of the bias and the precision of the approximate variance formulas we must refer to the literature. Since there may well be cases in which these approximations are too inaccurate we wish to put here on record a method which avoids these disadvantages under certain circumstances.

The situations in which our exact formulas will apply are characterized by the following conditions:

(a) The number of units in the j^{th} domain, jN , must be known.

(b) Primary sampling units (or, if the design is single stage, the sampling units) must have been drawn *with* replacements.

Two remarks concerning these conditions may be pertinent: Concerning (b), although there are many surveys in which primaries are actually drawn with replacement, this condition becomes of little importance when the sampling fraction of primaries is small as then there is practically no difference between drawing with — and drawing without — replacement. Small primary sampling fractions are very common in survey designs.

Concerning the very important condition (a), if the jN are known it is of course possible to obtain unbiased estimates of $j\bar{Y}$ by dividing the unbiased estimates of the jY (developed in the preceding sections) by the known jN . It is, however, well known from experience with such estimates that for most survey designs and populations they would have large variances. Confirmation of this may be sought by comparing for the examples quoted in the preceding sections $jN^{-2} \text{Var } j\hat{Y}$ with $\text{Var}(\tilde{y})$. Roughly speaking the former depends on the variation of jY — totals whilst the latter depends on the variation of jY — means.

Turning, therefore, again to the device of ratio estimation we will describe first the unbiased ratio estimators for population means translating them into estimators of domain means later.

Let us consider a survey design consisting of L strata ($h = 1, 2, \dots, L$) with population proportions P_h and let us assume that n_h primaries $t = 1, 2, \dots, n_h$ have been drawn with replace-

ment from the h^{th} stratum. Denote by y_{hti} a character + attached to the i^{th} unit in the t^{th} primary of the h^{th} stratum and by \hat{y}_{ht} the unbiased estimate of \bar{Y}_h from the t^{th} primary only. The functional form of \hat{y}_{ht} will depend on the design details.

Denote by

$$\bar{y}_h = n_h^{-1} \sum_{t=1}^{n_h} \hat{y}_{ht} \quad [54]$$

and by

$$\bar{y}_{st} = \sum_{h=1}^L P_h \bar{y}_h \quad [55]$$

the stratum means of the \hat{y}_{ht} and the unbiased estimator of \bar{Y} .

Denote further by \hat{x}_{ht}, \bar{x}_h and \bar{x}_{st} the corresponding means for a second characteristic x for which the population mean \bar{X} is known and for which we assume that either $x_{hti} \geq 0$ or $x_{hti} = y_{hti} = 0$.

Introduce now the ratio variates

$$\hat{r}_{ht} = \begin{cases} \hat{y}_{ht}/\hat{x}_{ht} & \text{if } \hat{x}_{ht} > 0 \\ r^* & \text{if } \hat{x}_{ht} = \hat{y}_{ht} = 0 \end{cases} \quad [56]$$

where r^* is a constant conveniently chosen as discussed below. Finally introduce for these ratios the means \bar{r}_h and \bar{r}_{st} as above. The survey design can now be thought of as generating in each stratum an infinite joint population of $\hat{y}_{ht}, \hat{x}_{ht}, \hat{r}_{ht}$ by the infinitely repeatable process of

- drawing a single primary in accordance with the design,
- drawing from this primary the prescribed number of secondaries, tertiaries . . . in accordance with the design,
- computing $\hat{y}_{ht}, \hat{x}_{ht}, \hat{r}_{ht}$ for the drawn sample,
- replacing all units.

Following the lines of Hartley-Ross (1954) and Goodman-Hartley (1956) we can now construct the following estimators:

The unbiased ratio estimate of \bar{Y}/\bar{X}

$$r' = \bar{r}_{st} + \bar{X}^{-1} \hat{c} \quad [67]$$

The unbiased estimate of \bar{Y} :

$$y' = \bar{X} \bar{r}_{st} + \hat{c} \quad [68]$$

where

$$\hat{c} = (\bar{y}_{st} - \bar{x}_{st} \bar{r}_{st}) + \sum_h P_h^2 (n_h - 1)^{-1} \{ \bar{y}_h - \bar{x}_h \bar{r}_h \} \quad [69]$$

is an unbiased estimate, computed from the stratified sample, of the population covariance of \hat{x}_{ht} and \hat{r}_{ht} defined by

$$C = \sum_h P_h E (\hat{x}_{ht} - \bar{X}) (\hat{r}_{ht} - \bar{R}) \quad [70]$$

in which E refers to the infinite populations described above.

The estimates r' and y' as given by [67] and [68] are unbiased no matter what value was chosen for the constant r^* . The choice of r^* will, however, affect the variance of these estimates. In the practical situations in which these estimators will be applied the chance of $\hat{x}_{ht} = \hat{y}_{ht} = 0$ to arise will be very small and hence any uncertainty in the choice of r^* will have a small effect. The best choice is to take r^* equal to the best estimate

* The notation applies directly to two stage designs, the index i denoting the secondary unit. In three or higher stage designs i must be replaced by a multiple subscript.

of \bar{R} available in advance of sampling and must not be altered subsequently.

In order to apply these results to the estimation of the j^{th} domain mean we must use for the numerator variable

$${}_j y_{hti} = \begin{cases} y_{hti} & \text{if the } hti \text{ unit is in the } j^{\text{th}} \text{ domain} \\ 0 & \text{otherwise} \end{cases}$$

and for the denominator the count variable

$${}_j u_{hti} = \begin{cases} 1 & \text{if the } hti \text{ unit is in the } j^{\text{th}} \text{ domain} \\ 0 & \text{otherwise} \end{cases}$$

and substitute these in [66], [67] and [69]. It will be seen that for all selfweighting designs the leading term of [57] \bar{r}_{st} will be of the form

$$\bar{r}_{st} = \sum_h P_h n_h^{-1} \sum_i {}_j \bar{y}_{hti}$$

where ${}_j \bar{y}_{hti}$ is defined, as before, as the y -mean of the units (if any) in the j^{th} domain of the t^{th} primary of the h^{th} stratum and by r^* if there are no such units. The condition mentioned above, that the chance of ${}_j \bar{y}_{hti} = {}_j u_{hti} = 0$ should be small means that the present method of estimation should only be used if there is a reasonable chance that all domains are represented in the sample from each primary.

The unbiased estimation of ${}_j \bar{Y}$ utilizing the knowledge of the domain mean ${}_j \bar{X}$ of a concomitant variable will be based on [68] when it should be noted that y' will estimate ${}_j N {}_j \bar{Y} / N$, that \bar{X} will be ${}_j N {}_j \bar{X} / N$ so that the estimator of ${}_j \bar{Y}$ will be

$$\frac{N}{{}_j N} y' = {}_j \bar{X} \bar{r}_{st} + \frac{N}{{}_j N} c$$

where c is computed from [69] using the variates ${}_j y_{hti}$, ${}_j x_{hti}$ and the ratios ${}_j r_{hti}$. In this case, therefore, both ${}_j \bar{X}$ and ${}_j N$ must be known.

The exact variance formulas for both r' and y' have recently been derived by Goodman-Hartley in the special case of simple sampling from one stratum ($L = 1$). Their results indicate that the variance of y' is of a similar order of magnitude (sometimes larger, sometimes smaller) than that of the combined ratio estimator. For the stratified case results are not as yet available.

8. Inferences based on the variance formulas.

Two questions concerning inferences arise:

The variance formulas of the preceding sections are based on finite population sampling and, in particular, have the property to become zero when sampling is 100%. The obvious

logic of this property is that when all units of the population are sampled all means of all domains in such a population are known without error so that all non-zero contrasts between domain means are "significant". It must not be forgotten, however, that such inferences can only apply to the particular finite population under investigation and are not of any wider significance. For instance, in the example of Table 3, had the population of Des Moines been sampled 100% definite statements about differences in the number of persons/household of different Income Groups could have been made "without error" but such inferences would then only apply to the City of Des Moines at the time of the survey. If data from such a survey are to be used for (say) comparing "Income Groups" in "Cities like Des Moines" in "times like the present" in a more general sense, the finite population sampling theory must not be used. Whether the survey sample can be (artificially) regarded as one drawn from a wider population (in space or time) is a matter requiring special investigation and the onus of such an investigation falls upon those wishing to draw such wide inferences.

A second question concerns inferences to be drawn for the domains of the particular finite population for which the survey was planned:

In the preceding sections we have given estimates of domain means, their variances and estimates of these variances. We have also given variance formulas and estimates for the differences between estimated domain means. The question arises as to how these variance estimates are to be used for inferences to be drawn from the data. Even if it is not intended to carry out "tests of significance" there would remain the question of the computation of confidence intervals for the domain means and their contrasts. This question is, of course, not confined to the present issues but is an ever present difficulty in the theory of sample surveys. The customary procedure here is to employ normal theory approximations and to appeal to the respective central limit theorems for finite populations and to the relatively large sample sizes which are available in these situations. The results which are available on these problems are very restricted and require development and it is clearly not possible in this context to dwell upon this issue of much wider impact.

REFERENCES

- COCHRAN, W. G., (1953). — *Sampling Techniques*. John Wiley and Sons, Inc. New York.
- GOODMAN, L., HARTLEY, H. O., (1956). (Unpublished manuscript).
- HARTLEY, H. O., ROSS, A. — *Unbiased Ratio Estimators*, 1954. *Nature*, 174, 270.
- KISH, L. and HESS, IRENE, (1955). — *On Variances of Ratios and Their Differences in Multistage Sampling*. (Paper presented at the New York Meeting of A.S.A. Dec. 1955).
- STEPHAN, F. F., (1941). — *Stratification in Representative Sampling*. *J. of Marketing* (1941), 38.
- YATES, F., 1949 (2nd edition 1953). — *Sampling Methods for Censuses*. Surveys, London, Charles Griffin & Co.

ACKNOWLEDGMENT

This paper has appeared as a contribution to a volume published by the Institute of Statistics in honour of Professor Corrado Gini, upon the occasion of his retirement from the Chair of Statistics at the University of Rome. It is reproduced here through the kind permission of the Institute of Statistics and of its Director, Professor Vittorio Castellano. Professor Hartley presented a summary of this paper at the Annual Meeting.