Robert D. Tortora, USDA

Nicholas J. Ciancio, USDA

Introduction

Recently work has been done comparing various approximate optimum stratification techniques. Hess, et.al., [6] and Cochran, [2] compared various techniques for actual populations. In [2], Cochran was concerned with eight different populations ranging from income per tax return, population of US cities, resources of commercial banks, number of farms per area sampling unit, and proportion of gross bank loans. In [6], the stratification and primary estimation variable was size of hospital. A brief discussion examines other estimation variables with high correlations (> .9) with the stratification variable, Kish, et.al., [1] compared various stratification techniques for a specified bivariate population where the stratification is carried out on an auxiliary variable X and estimation is made for a variable Y. Kpedeko, [7] in a review of the literature on stratification techniques calls for further empirical studies to evaluate some of these methods for different types of data.

This paper compares five approximate optimum stratification techniques when an auxiliary variable is used for stratification and when one is interested in estimating crop acreage and livestock totals. The stratification techniques are 1) cum \sqrt{f} , 2) Durbin, 3) Ekman, 4) Sethi, and 5) Equal Aggregate Output (EAO). The Statistical Reporting Services' (SRS) area frame is used in two States to make the comparisons. In the area frame the land area is classified (stratified) according to land use in order to achieve homogeneity within strata. The sampling unit is a segment, which is a piece of land with boundaries delineated on a map. Every parcel of land within a segment is accounted for during a survey.

The stratification variable for the area frame is the percent of land under cultivation. For each segment this is defined as the total cropland in acres in the segment divided by the total acres in the segment times 100. The crops acreage variables of interest are the three most important income-producing crops for US farmers, vis., corn, wheat, and soybeans. Similarly, the important livestock variables are cattle and hogs and these variables will also be studied.

The data used in this study are from the 1975 June Enumerative Survey for Ohio and Kansas. The segments are from the agricultural strata and the population sizes are N=252 and N=435 segments in Ohio and Kansas, respectively. Even though the above five commodities are the most important within Ohio and Kansas, these States differ demographically and geographically. Kansas is a more homogeneous farm state with more area under irrigation. The average size of farm in Kansas is larger (616 acres vs. 150 acres). Ohio has more farms (117,000 vs. 81,000) and less land in farms $(17.5 \times 10^6 \text{ acres vs } 50 \times 10^6 \text{ acres})$. The segment size is Kansas ranges from 1 to 4 square miles while in Ohio the segment size is $\frac{1}{2}$ to 1 square miles.

Optimum allocation for fixed sample size is used to determine sample sizes in the strata. This technique is the one used by SRS.

The comparisons are made for 2, 3, 4 or 5 strata. Currently SRS uses four strata with stratum boundary values 15%, 50%, and 75%. The variances for the approximate techniques will be compared to the variance under the current SRS technique. The total number of strata is held to 5 for two reasons. First, to strain the stratification techniques, which depend more or less on the assumption that the number of strata L is reasonably large, so that within a stratum the frequency function can be assumed to be rectangular. Secondly, for practical purposes of frame construction, it is very difficult to efficiently divide the area frame into a large number of strata.

The Approximate Methods

Let X_0, X_1, \ldots, X_L be the stratum boundaries, the strata being numbered 1, 2, ..., L. let S_h be the standard deviation is stratum h and W_h = N_h/N be the ratio of the number of sampling units in stratum h to the total number in the population. The usual estimate of the population total is

$$y_{st} = N\Sigma W_h \overline{y}_h$$

where \overline{y}_{h} is the sample mean in statum h. Its variance is

$$\mathbb{V}(\mathbf{y}_{st}) = \mathbb{N}^{2} \mathbb{E} \mathbb{W}_{h}^{2} S_{h}^{2} \quad (\frac{1}{n_{h}} - \frac{1}{N_{h}}).$$

For a fixed total sample size of n, $V(y_{st})$ is minimized by taking $n_n = nN_h S_h/\Sigma N_h S_h$. The minimum variance is

(1)
$$V_{\min}(y_{st}) = \frac{1}{n} (\Sigma N_h S_h)^2$$
,

ignoring the fpc. Equation (1) becomes the basis for further caculations.

A discussion of the actual implementation of the approximate stratification techniques can be found in [2], [6] or [7].

The Study Variables

Table 1 gives the shapes of the stratification variables in the two States. It is simply the percentage of the total number of segments lying within each tenth of the range of the percent of land under cultivation. Table 1: Percentage of segments falling into successive tenths of the stratification variable.

	Ohio	Kansas
0-10	8.73	10.11
10-20	7.14	3.91
20-30	5.56	5,06
30-40	5.16	4.60
40-50	9.52	7.82
50-60	7.52	8.74
60-70	9.52	10.57
70-80	14.68	12.64
80-90	15.87	15.40
90-100	16.28	21.15

The distribution in Ohio is closest to a rectangular distribution, while in Kansas the distribution is close to being two-tailed. Cum \sqrt{f} , Durbin and Ekman compute approximately the same stratum boundary values, as do EAO and Sethi. Also all techniques seem to be relatively insensitive to the distribution of the stratification variable although in Kansas the values are higher than in Ohio.

Table 2 and 3 give the shape of the frequency distribution of the variables to be estimated.

Table 2: Percentage of segments having crop acreage falling into given classes by State, (Ohio/Kansas).

acres	corn	wheat	soybeans
0-10	17.46/76.78	34.92/13.79	31.75/80.46
10-50	26.59/ 8.04	38.10/ 9.66	19.05/ 8.50
50-100	36.11/ 5/98	22.62/12.18	24.60/ 6.21
100-250	19.44/ 4.83	4.36/30/11	23.81/ 4.60
> 250	0.00/ 4.37	0.00/34.26	0.79/ 0.03

Table 3: Percentage of segments having livestock numbers falling into given classes by States, (Ohio/Kansas).

Number of	livestock	cattle	hogs
0		69.65/85.29	28.97/31.49
1 -	50	21.43/ 8.28	40.48/25.98
51 -	100	3.97/ 2.53	18.65/22.30
101 -	500	4.36/ 3.45	11.90/19.54
>	500	1.59/ 0.45	0.00/ 0.69

The tables show that in Ohio corn and wheat appear unimodal while soybeans appear bimodal. In Kansas corn and soybeans are skewed to the left while wheat is skewed to the right. The general shape of the distribution for hogs are the same. For cattle, the distributions are skewed to the left with the distribution in Ohio being fatter.

Finally Table 4 lists the estimated correlation coefficients between the stratification variable and the variables of interest. Note that none of the correlations are near the correlations in [6] and thus should strain the techniques. Table 4: Estimated correlation coefficients between the stratification variable and the variables of interest, (Ohio/Kansas).

	corn	wheat	soybeans
Percent of land under cultivation	<u>.554</u> .240	<u>.614</u> .609	<u>.652</u> .028
	cattle	hogs	
percent of land under cultivation	.118	$\frac{216}{356}$	

The correlation coefficients between the crops and the auxiliary variable are consistent in Ohio while in Kansas wheat has the highest correlation. In both states the magnitude of the correlations between livestock and percent cultivated are small.

Comparison of the Rules

Equation (1) is used as the basis for the comparisons for each technique. The variance under the approximate stratification technique is compared to the variance obtained under the present stratification used by SRS. The results are presented in Table 5.

The separation of boundary values (cum \sqrt{f} , Durbin, and Ekman vs. EAO and Sethi) is reflected in Table 5. As the number of strata increases the differences between the techniques does not change as much for the crop variables, i.e., for those with higher correlation coefficient with the stratification variable and where the stratification variable is more nearly rectangular (Ohio). Whenever EAO or Sethi have a smaller ratio for 2 strata than either $\operatorname{cum}\sqrt{f}$, Durbin, or Ekman they retain their smaller ratio as the number of strata increases. For the negatively correlated variables in both states $\operatorname{cum}\sqrt{f}$, Durbin and Ekman perform much better than the other techniques. For the highest positively correlated variables (soybeans in Ohio and wheat in Kansas) $\operatorname{cum}\sqrt{f}$, Durbin, and Ekman perform well. To get a feel for the performance of the techniques across all variables of interest a plot of the technique(s) (with smallest ratios from Table 5 vs. the correlation coefficients (r) is given in Figure 1. For the range of correlation coefficients we see that $cum\sqrt{f}$, Durbin and Ekman perform well for negatively correlated variables as well as for moderately correlated values (r > .5). For small positive values of r all techniques seem to perform on a par.

Table 5: Variance under the approximate stratification technique divided by the variance under the current SRS stratification (Ohio/Kansas).

.

technique	strata	corn	wheat	soybeans	cattle	hogs
cum√f	2	$\frac{3.434}{2.499}$	$\frac{3.229}{2.963}$	$\frac{3.165}{3.012}$	$\frac{3.224}{2.341}$	$\frac{3.528}{3.374}$
	3	$\frac{1.334}{1.133}$	$\tfrac{1.221}{1.316}$	$\tfrac{1.261}{1.282}$	$\frac{1.090}{1.149}$	$\frac{1.546}{1.470}$
	4	$\frac{0.781}{0.643}$	$\tfrac{0.676}{0.719}$	$\frac{0.699}{0.703}$	$\frac{0.707}{0.603}$	$\frac{0.882}{0.658}$
	5	$\frac{0.471}{0.394}$	$\frac{0.434}{0.460}$	$\frac{0.421}{0.442}$	$\frac{0.437}{0.342}$	$\frac{0.773}{0.461}$

TADYE 2. (OAn c)	Table	5:	(Con'	t)
------------------	-------	----	-------	----

technique strata corn wheat soybeans cattle hogs

Durbin	2	$\frac{3.595}{2.499} \frac{3.959}{2.963}$	$\frac{3.334}{3.012}$	$\frac{3.239}{2.341}$	$\frac{3.514}{3.374}$
	3	$\frac{1.440}{1.046} \frac{1.357}{1.236}$	$\tfrac{1.418}{1.266}$	$\frac{1.260}{0.989}$	$\tfrac{1.622}{1.398}$
	4	$\frac{0.806}{0.643} \frac{0.724}{0.719}$	$\frac{0.721}{0.703}$	$\frac{0.718}{0.603}$	$\tfrac{0.861}{0.658}$
	5	$\frac{0.483}{0.377} \frac{0.452}{0.439}$	$\frac{0.433}{0.420}$	$\frac{0.446}{0.285}$	0.538 0.444
Ekman	2	$\frac{3.909}{2.499} \frac{3.229}{2.963}$	$\frac{3.932}{3.012}$	$\frac{3.024}{2.341}$	$\frac{3.527}{3.374}$
	3	$\frac{1.369}{1.133} \frac{1.254}{1.316}$	$\tfrac{1.297}{1.282}$	$\frac{1.210}{1.149}$	$\frac{1.566}{1.470}$
	4	$\frac{0.753}{0.714} \frac{0.656}{0.798}$	$\frac{0.638}{0.750}$	$\frac{0.650}{0.720}$	$\frac{0.865}{0.889}$
	5	$\frac{0.498}{0.440} \frac{0.464}{0.481}$	$\frac{0.425}{0.428}$	$\frac{0.467}{0.424}$	$\frac{0.554}{0.429}$
EAO	2	$\frac{4.281}{2.316} \frac{4.055}{3.776}$	$\frac{3.981}{3.253}$	$\frac{3.139}{1.601}$	$\frac{5.187}{4.783}$
	3	$\frac{1.955}{0.818} \frac{1.639}{1.604}$	$\tfrac{1.642}{1.451}$	$\frac{1.541}{0.803}$	$\frac{2.494}{2.404}$
	4	$\frac{1.045}{0.450} \frac{0.942}{0.912}$	$\frac{0.934}{0.836}$	$\frac{0.844}{0.343}$	$\tfrac{1.650}{1.433}$
	5	$\frac{0.723}{0.250} \frac{0.548}{0.632}$	<u>0.497</u> 0.522	$\tfrac{0.560}{0.200}$	$\tfrac{1.099}{1.034}$
Sethi	2	$\frac{4.281}{1.963} \frac{4.055}{3.320}$	$\frac{3.981}{3.480}$	$\frac{3.139}{1.635}$	$\frac{5.013}{4.396}$
	3	$\frac{1.915}{0.896} \frac{1.610}{1.556}$	$\tfrac{1.439}{1.698}$	$\frac{1.511}{0.783}$	$\frac{2.425}{2.165}$
	4	$\frac{1.048}{0.497} \frac{0.925}{0.847}$	$\frac{0.945}{0.937}$	$\tfrac{0.842}{0.339}$	$\frac{1.696}{1.284}$
	5	$\frac{0.758}{0.260} \frac{0.590}{0.604}$	<u>0.529</u> 0.583	$\frac{0.648}{0.195}$	$\frac{1.230}{0.936}$

From Figure 1 we see $\operatorname{cum}\sqrt{f}$ performs best 18 times, but Durbin is best 14 times, and Ekman is best 9 times. $\operatorname{Cum}\sqrt{f}$ peforms well over the range of r, Durbin does well with smaller values of r, and Ekman does well with larger values of r. Figure 2 presents a graphic display of the worst of the best for the regions were the three techniques $\operatorname{cum}\sqrt{f}$, Durbin, and Ekman perform well. Figure 2 presents the technique with the largest ratio from Table 6. The trends exhibited here indicate that $\operatorname{cum}\sqrt{f}$ can give larger variances for negative r, Durbin for moderate r, and Ekman across the range of r.

Dalenius [3] suggested the approximation $V_L/V_{L-1} = (L-1)^2/L^2$ (for rectangular distributions) to quantify the gains caused by stratification. For L = 2, 3, 4, and 5 we get from the formula 0.250, 0.444, 0.562 and 0.640 respectivly. Table 6 presents the average gain by crop or livestock for each technique.

In general, the average gain is slightly more than that estimated by Dalenius. There is less gain in precision for the variables than with lower correlation (livestock in Ohio) with the stratification variable than with those with higher correlation (crops in Ohio). Comparing gains against the distribution of the stratifying variable we see that there are more gains (25 vs 15) for the unimodal distribution (Ohio). Defining any gain exceeding $(L-1)^2/L^2$ as significant, we see that for the unimodal distribution there are more significant gains (22 vs. 14). Finally, examining significant gains by correlation and by technique we see that cum/f does best for the higher correlations (crop in Ohio and Kansas), Durbin performs about the same for crop and livestock and the remaining three produce more significant gains for the lower correlation (livestock in Ohio and Kansas).

Table 6: The average gain $\rm V_L/V_{L-1}$ for crops and livestock by straitification technique (Ohio/Kansas)

		crops	livestock
cum f	2	<u>.188</u> .231	.248
	3	<u>.388</u> .440	<u>.388</u> .464
	4	<u>.546</u> .554	<u>.610</u> .486
	5	<u>.616</u> .627	<u>.747</u> .634
Durbin	2	.202	<u>.248</u> .220
	3	<u>.401</u> .419	<u>.425</u> .418
	4	.54 1 .584	<u>.550</u> .540
	5	<u>.608</u> .598	<u>.623</u> .574
Ekman	2	<u>.176</u> .231	<u>.287</u> .220
	3	<u>.470</u> .441	.422
	4	<u>.522</u> .607	.545
	5	<u>.679</u> .597	<u>.679</u> .536
EAO	2	<u>.263</u> .256	<u>.158</u> .226
	3	.424	.496
	4	<u>.559</u> .508	<u>.606</u> .512
	5	<u>.602</u>	<u>.662</u>

nu

		crops	livestock
Sethi	2	.236	<u>.292</u> .216
	3	<u>.402</u> .457	<u>.482</u> .486
	4	<u>.593</u> .558	<u>.628</u> .512
	5	<u>.640</u> .620	<u>.747</u> .652

Summary

This study compared five approximate techniques for stratification in an agricultural setting. The comparisons were based on area sampling units from two States. The stratification variable (percent of land under cultivation) was different from the variables to be estimated (corn, wheat, soybeans, cattle and hogs).

The rules divided themselves into two groups based on stratum boundary values, $cum\sqrt{f}$, Durbin, Ekman, and Equal Aggregate Output, Sethi. Comparisons were based on variances obtained using the current SRS stratification. $Cum\sqrt{f}$, Durbin, and Ekman performed well for variables either with negative correlations or moderate positive correlations with the stratification variable. All five rules were comparable for small positive correlations.

Using Dalenius' approximation, $(L-1)^2/L^2$, for gains due to increasing the number of strata it was found that the most significant gains were produced when the stratification variable was unimodal. Ekman, Equal Aggregate Output, and Sethi had more significant gains for variables not highly correlated with the stratification variable, gains and cum/f produced more significant gaims with the higher correlated variables. It was found that the approximation $(L-1)^2/L^2$ was an overestimate of the gains due to increasing the number of strata (concurring with the results in [2]).

References

- [1] Anderson, D.W., Kish, L., Cornell, R.G.
 "Quantifying Gains From Stratification for Optimum and Approximately Optimum Strata Using a Bivariate Normal Model", <u>Tech. Report</u> 4, Department of Biostatistics, The University of Michigan, (1975), Ann Arbor, Michigan 48104
- [2] Cochran, W.G. "Comparison of Methods for Determining Stratum Boundaries", <u>Bulletin</u> of the International Statistical Institute, 38(2), Tokyo (1961), pp. 345-358.
- [3] Dalenius, T. <u>Sampling in Sweden</u>, Chapter 8, Almquist and Wiksell, Stockholm (1957).
- [4] Dalenius, T. and Hodges, J.L., Jr. "Minimum Variance Stratification", <u>Journal of the</u> <u>American Statistical Association</u>, 54 (1959), pp. 88-101
- [5] Hansen, M.H., Hurwitz, W.N., Madow, W.G. <u>Sample Survey Methods and Theory</u>, John Wiley and Sons, Inc. 1953.
- [6] Hess, I., Sethi, V.K., and Balakrishnan, T.R. "Stratification: A Practical Investigaion", <u>Journal of the American Statistical Associa-</u> <u>tion, 61, (1966) pp. 74-90.</u>
- [7] Kpedekpo, G.M.K. "Recent Advances on Some Aspects of Stratified Sample Design. A Review of the Literature", <u>Metrika</u>, 20 (1073), pp. 54-64.
- [8] Mahalanobis, P.C. "Some Aspects of the Design of Sample Surveys. <u>Sankhya</u> 12 (1952) pp. 1-7.

	356	216	.028	.077	.118	.240	.554	.609	.614	.652	r
	2CDE	D	CDE	EAO	EAOS	S	С	CDE	CE	С	
	3D	С	D	EAO	С	EAO	С	D	С	С	
strata	4 - CD	D	CD	S	E	EAO	Ε	CD	Е	Е	
mber of	5 - E	D	D	S	С	EAO	С	D	С	С	

Figure 1: Best technique (smallest ratio from Table 6) vs. Correlation coefficient (ordered by increasing magnitude). $C=cum\sqrt{f}$, D=Durbin, E+Ekman, EAO=EAO, S=Sethi.

number of strata	5 +C 4 +E 3 +CE	C C D	E D D	E E CE	E D D	ם ס ס	
	2 CDE	С	E	CDE	D	E	
	356	216	.554	.609	.614	.652	

-.356 -.216 .554 .609 .614 .652 r Figure 2: Worst technique (largest ratio from Table 6) vs. correlation coefficient over ranges where cum√f, Durbin and Ekman perform well.