SOME SMALL SAMPLE RESULTS FOR THE VARIANCE OF A RATIO

Elizabeth Lauh and Wm. H. Williams Bell Telephone Laboratories, Incorporated

1. Introduction

Ratio estimators are used to increase the reliability of the sample estimates. For each sampled unit an auxiliary variate, x, is observed in addition to the variate of interest, y, so that for each unit a pair (y,x) is obtained. From a random sample of n pairs (y_1,x_1) i = 1,2,...n a common survey problem is the estimation of population mean, \overline{Y} , subject to the assumption that the population mean \overline{X} is known exactly. Ratio estimators are designed to utilize this information and the most common ratio estimator of \overline{Y} is $\widetilde{y} = (\overline{y}/\overline{x})\overline{X}$, which is generally biased. A study of certain statistical properties of the ratio-of-sample-means estimator of \overline{Y} is the main interest of this paper.

The sampling variance of \widetilde{y} depends upon the sampling variance of y/x for which there is no known exact expression. An approximation to this sampling variance, obtained by the use of a Taylor expansion, is well-known and can be found in most textbooks. A sample variance is usually formed by substituting sample expressions for the population quantities which appear in the sampling variance approximation. For large samples, the relationship between this sampling variance approximation and the sample variance thus obtained has been examined analytically and a relatively clear picture of the situation obtained. For small samples, the relationship is not so clear; nor is the relationship between the sampling variance approximation and the exact sampling variance. Unfortunately, these questions do not lend themselves to simple and productive analytic study and so this paper presents some Monte Carlo results for the sample sizes n = 2(1)9.

2. The Ratio Estimator Approximations 2.1 The Bias

Approximations to both the bias and sampling variance of y/x can be found by the use of Taylor expansions. If the negative exponent of the identity

$$\frac{\overline{y}}{\overline{x}} = \frac{\overline{Y}}{\overline{X}} \left(1 + \frac{d\overline{y}}{\overline{Y}} \right) \left(1 + \frac{d\overline{x}}{\overline{X}} \right)^{-1}, \qquad (1)$$

where

$$d\bar{y} = \bar{y} - \bar{Y}$$
 and $d\bar{x} = \bar{x} - \bar{X}$,

is expanded as

$$\left(1 + \frac{d\bar{x}}{\bar{x}}\right)^{-1} = \sum_{k=0}^{\infty} (-1)^{k} \left(\frac{d\bar{x}}{\bar{x}}\right)^{k}$$
(2)

then a first approximation to the bias is

$$\operatorname{Bias}_{1}\left(\frac{\bar{y}}{\bar{x}}\right) = \frac{\bar{Y}}{\bar{x}} \left(C\bar{x}\bar{x} - C\bar{y}\bar{x} \right), \qquad (3)$$

where

and

 $C\bar{y}\bar{x}$ denotes $Cov(\bar{y},\bar{x})/\bar{Y}\bar{X}$

$$C\bar{x}\bar{x}$$
 denotes $V(\bar{x})/\bar{X}^2$.

The expansion of Equation (2) is valid for $\left|\frac{d\bar{x}}{\bar{x}}\right| < 1$ and terms up to the second order have been retained to obtain Equation (3).

2.2 The Variance

The variance of $\overline{y}/\overline{x}$ is, by definition,

$$V(\bar{y}/\bar{x}) = E(\bar{y}/\bar{x})^2 - [E(\bar{y}/\bar{x})]^2, \quad (4)$$

and it is convenient to consider the two terms on the right-hand side separately. Squaring both sides of Equation (1) and using the expansion,

$$\left(1 + \frac{d\bar{x}}{\bar{x}}\right)^{-2} = \sum_{k=0}^{\infty} (-1)^{k} (k+1) \left(\frac{d\bar{x}}{\bar{x}}\right)^{k}, \quad (5)$$

the second order approximation to $E(\bar{y}/\bar{x})^2$ is obtained as,

$$E\left(\frac{\bar{y}}{\bar{x}}\right)^2 = \frac{\bar{y}^2}{\bar{x}^2} \left(1 + 3C\bar{x}\bar{x} - 4C\bar{y}\bar{x} + C\bar{y}\bar{y}\right).$$
(6)

Next, squaring $E\left(\frac{\bar{y}}{\bar{x}}\right)$, expanding, and dropping third and higher order terms gives $\left[E(\bar{y}/\bar{x})\right]^2 = \frac{\bar{y}^2}{\bar{x}^2} (1 + 2C\bar{x}\bar{x} - 2C\bar{y}\bar{x}),$ (7)

so that subtraction from Equation (6) yields a first approximation to V(y/x) as

$$V_1(\bar{y}/\bar{x}) = \frac{\bar{y}^2}{\bar{x}^2} (C\bar{y}\bar{y} + C\bar{x}\bar{x} - 2C\bar{y}\bar{x}).$$
 (8)

The large sample variance of $\widetilde{\mathbf{y}}$ can then be written as,

$$V_{1}(\tilde{y}) = \frac{1}{n} \bar{Y}^{2}(Cyy + Cxx - 2Cyx).$$
(9)

These large sample approximations to the bias and variance of \tilde{y} are those usually found in the literature, [1,3], and seem to hold for large n, see for example, Cochran, [1, p. 114].

An approximation to the mean square error can be found by the same procedure, see Kish [3] and Sukhatme [5].

The approximation \boldsymbol{V}_1 can be written as,

$$V_{1}(\tilde{y}) = \frac{1}{n} V(y_{1} - Qx_{1}), \qquad (10)$$

where

$$Q = \overline{Y}/\overline{X}$$
.

This suggests the sample estimate,

$$v_1 = \frac{1}{n(n-1)} \sum_{i=1}^{n} (y_i - qx_i)^2$$
, (11)

with

$$q = \bar{y}/\bar{x}$$
.

This sample estimator $v_1(\tilde{y})$ has a bias of order $O(n^{-1})$, [1, p. 119].

3. Small Sample Estimation

When the sample size is large, the approximations described in Section 2 have good accuracy, and it is not necessary to bring in the higher moments by inclusion of more terms in the expansions. With small samples however, the approximations do not always have a known accuracy and it is natural to consider improving the estimates by including some higher order terms.

By expanding Equation (2) up to fourth order terms, a second approximation to the bias of \tilde{y} is obtained as,

$$\operatorname{Bias}_{2}(\widetilde{\mathbf{y}}) = \operatorname{Bias}_{1}(\widetilde{\mathbf{y}}) + \overline{\mathbf{y}} \left[-\frac{\mu_{03}}{\overline{\mathbf{x}}^{3}} + \frac{\mu_{04}}{\overline{\mathbf{x}}^{4}} + \frac{\mu_{12}}{\overline{\mathbf{y}}\overline{\mathbf{x}}^{2}} - \frac{\mu_{13}}{\overline{\mathbf{y}}\overline{\mathbf{x}}^{3}} \right]$$
(12)

where $\mu_{rs} = E(\bar{y}-\bar{Y})^r(\bar{x}-\bar{X})^s$.

A second approximation to the variance can be found by expanding Equation (5) up to fourth order terms and again treating each term of Equation (4) separately. This gives the following equation as a second approximation to $V(\tilde{y})$.

$$v_{2}(\tilde{y}) = v_{1}(\tilde{y}) + \bar{y}^{2} \left[-\frac{2\mu_{03}}{\bar{x}^{3}} + \frac{3\mu_{04}}{\bar{x}^{4}} - \frac{\mu_{02}^{2}}{\bar{x}^{4}} + \frac{\mu_{02}}{\bar{x}^{4}} + \frac{\mu_{02}}{\bar{x}^{2}} + \frac{3\mu_{11}}{\bar{y}\bar{x}^{2}} - \frac{6\mu_{13}}{\bar{y}\bar{x}^{3}} + \frac{2\mu_{02}\mu_{11}}{\bar{y}\bar{x}^{3}} - \frac{2\mu_{21}}{\bar{y}^{2}\bar{x}} + \frac{3\mu_{22}}{\bar{y}^{2}\bar{x}^{2}} - \frac{\mu_{11}^{2}}{\bar{y}^{2}\bar{x}^{2}} \right]$$

$$(13)$$

In the Monte Carlo results presented in Section 4, y is generated as a linear function of x and so the consideration of some of the results when y is a linear function of x, is of interest. In this case, y = A+Bx, Equation (12) becomes

$$\operatorname{Bias}_{2}(\tilde{\mathfrak{Y}}) = \frac{A}{\bar{\mathfrak{X}}} \left(\frac{\mu_{2}}{\bar{\mathfrak{X}}} - \frac{\mu_{3}}{\bar{\mathfrak{X}}^{2}} + \frac{\mu_{4}}{\bar{\mathfrak{X}}^{3}} \right), \quad (14)$$

where the moments of \bar{x} about \bar{X} are denoted with a single subscript, i.e.

 $\mu_{rs} = B^r \mu_{r+s}$. In terms of the moments of x_1 about the mean,

$$Bias_{2}(\tilde{y}) = \frac{1}{n} \left[\frac{A\mu_{2}(x)}{\bar{x}^{2}} \right] + \frac{1}{n^{2}} \frac{A}{\bar{x}} \left[\frac{-\mu_{3}(x)}{\bar{x}^{2}} + \frac{\mu_{4}(x)}{n\bar{x}^{3}} + \frac{3(n-1)\mu_{2}^{2}(x)}{n\bar{x}^{3}} \right],$$

(15)

where on the first term on the right is $Bias_1(\tilde{y})$ and the second term is $O(n^{-2})$.

In this case of a linear relationship between y and x, the difference between the two approximations to the variance is,

which in terms of the moments of \mathbf{x}_{1} about the mean is

$$V_{2-1}(\tilde{y}) = \frac{1}{n^2} \frac{A^2}{\bar{x}^2} \left(-\frac{2\mu_3(x)}{\bar{x}} + \frac{3\mu_4(x)}{n\bar{x}^2} + \left(\frac{8n-9}{n}\right) \frac{\mu_2^2(x)}{n} \right)$$
(17)

This difference is positive in symmetric populations and is also $O(n^{-2})$.

It should be noted that both the bias and the variance are direct functions of the y intercept, A, so that if A = 0, \tilde{y} is exactly unbiased and $V_1(\tilde{y}) = V_2(\tilde{y}) = 0$. This is true for any sample size. If y is a nonlinear function of x, say quadratic, $y_1 = A + Bx_1 + Cx_1^2$, then the bias and variance of \tilde{y} are affected not only by the y-intercept A, but also by the nonlinear coefficient, C. Then, even if A = 0, neither the bias nor the variance vanishes.

A higher order approximation to the mean square error can be obtained by combining the identity,

$$\left(\frac{\overline{y}}{\overline{x}} - \frac{\overline{y}}{\overline{x}}\right)^2 = \left(\frac{d\overline{y} - Qd\overline{x}}{\overline{x}}\right)^2 \left(1 + \frac{d\overline{x}}{\overline{x}}\right)^{-2},$$

with the expansion of Equation (5) and including terms up to the fourth order. This gives,

$$MSE_{2}(\tilde{y}) = \tilde{x}^{2} E(\tilde{y}/\tilde{x}-Q)^{2}$$

$$= \tilde{y}^{2} \left(\frac{\mu_{20}}{\tilde{y}^{2}} - \frac{2\mu_{11}}{\tilde{y}\tilde{x}} + \frac{\mu_{02}}{\tilde{x}^{2}} - \frac{2\mu_{03}}{\tilde{x}^{3}} + \frac{3\mu_{04}}{\tilde{x}^{4}} + \frac{\mu_{12}}{\tilde{y}\tilde{x}^{2}} - \frac{6\mu_{13}}{\tilde{y}\tilde{x}^{3}} - \frac{2\mu_{21}}{\tilde{y}^{2}\tilde{x}} + \frac{3\mu_{22}}{\tilde{y}^{2}\tilde{x}^{2}}\right), \quad (18)$$

which is an expression that appears in Kish [3] and Sukhatme [5]. The difference between $V_2(\tilde{y})$ and $MSE_2(\tilde{y})$ is

$$\mathbf{V}_{2}(\mathbf{\tilde{y}}) - \mathrm{MSE}_{2}(\mathbf{\tilde{y}}) = -\frac{\mathbf{\tilde{Y}}^{2}}{\mathbf{\tilde{X}}^{2}} \left(\frac{\mu_{02}}{\mathbf{\tilde{X}}} - \frac{\mu_{11}}{\mathbf{\tilde{Y}}}\right)^{2}$$
 (19)

which is the square of the estimate of bias obtained in Equation (3), and, since $\mu_{11} = \rho \sigma_{\overline{x}} \sigma_{\overline{y}}$,

$$V_{2}(\tilde{y}) - MSE_{2}(\tilde{y}) = -\bar{Y}^{2}C_{\bar{x}}^{2}(C_{\bar{x}} - \rho C_{\bar{y}})^{2} \quad (20)$$

which is always negative and becomes negligible for large n.

Finally, when y and x are bivariate normal, $MSE_2(\breve{y})$ reduces to,

$$V_1(\tilde{y}) \cdot (1+3C\bar{x}\bar{x}) + 6 \operatorname{Bias}_1^2(\tilde{y}).$$
 (21)

The sample estimate $v_2(\tilde{y})$ of $V_2(\tilde{y})$ is obtained by estimating each term of Equation (13) from the sample. Similarly, an estimate $mse(\tilde{y})$ of $MSE(\tilde{y})$ can be found from Equation (18), or, if y and x are bivariate normal, from Equation (21).

Recently, estimators which reduce bias have been constructed by splitting up samples, see Quenouille [4]. Tukey [6] later discussed the split-sample procedure from the point of view of ease in variance computations. The method consists of constructing estimates $q_{(1)}$, i = 1,2,...n, based on all but the ith observation and then forming

$$q_{(1)}^* = nq - \overline{n-1} q_{(1)},$$
 (22)

where ${\bf q}$ is the estimator based on all the observations. The estimator

$$q^* = \sum_{i=1}^{n} q^*_{(i)} / n,$$
 (23)

then has a sample variance given approximately by,

$$v(q^*) = \frac{1}{n(n-1)} \sum_{i=1}^{n} (q^*_{(i)} - q^*)^2.$$
 (24)

Thus in the case of the ratio-of-means estimator $% \left({{{\left({{{{{\bf{n}}}} \right)}_{{{\bf{n}}}}}}} \right)$

$$q_{(1)} = \frac{\sum_{j=1}^{n} y_{j} - y_{j}}{\sum_{j=1}^{n} x_{j} - x_{j}},$$
 (25)

 $q = \overline{y}/\overline{x}$, and an estimator of \overline{Y} is given by $\widetilde{y}^* = q^*\overline{x}$, with an approximate sample variance given by,

$$\mathbf{v}(\mathbf{\tilde{y}^{*}}) = \frac{1}{n(n-1)} \sum_{i=1}^{n} (q^{*}_{(i)} - q^{*})^{2}.$$
 (26)

4. Monte Carlo Study

To examine the problem of variance estimation for small sample sizes, random samples of size n = 2(1)9 were drawn from a population of fifty thousand pairs (y,x). The ratio $\overline{y}/\overline{x}$ was used as an estimator of $\overline{Y}/\overline{X}$ but the results can be applied directly to \widetilde{y} .

It is a very large undertaking to draw all possible samples for a given sample size so instead one thousand samples were drawn from the population and $q = \overline{y}/\overline{x}$ computed for each. The variance of the one thousand q's should be a close approximation to the true sampling variance. The split-sample estimator q^* was found as well as q, the four variance estimators $v_1(q)$, $v_2(q)$, $mse_2(q)$, and $v(q^*)$ were also calculated for each sample.

A population of fifty thousand pairs, (y_1, x_1), was constructed with $x_1 N(10, 4)$, and y_1 defined as $5(x_1+e_1)$ where e_1 is a random error distributed as standard normal independently of x_1 , y is therefore N(50,125), and the correlation between y and x is .89.

Normal probability plots were made of the one thousand sample estimates of q and q*. The first three figures show the plots of q for sample sizes 3, 6 and 9. The figure numbers in all the plots correspond to sample sizes. These plots show that even for small sample sizes the ratio \bar{y}/\bar{x} is essentially normal when the original y and x populations are normal. The indications of the probability plots of q* were the same as those for q.

Table 1 gives the mean and variance of the sampling distributions of q and q* for the different sample sizes. The expected values are essentially the same, but the precision of q seems to increase over that of q* as n decreases.

TABLE 1

Mean and Variance of 1000 Samples of q and q* for Small Sample Size n

		q	q *		
n	Mean	Variance	Mean	Variance	
23456789	4.9975 4.9889 4.9971 4.9928 5.0055 4.9932 5.0041 5.0026	.1235 .0896 .0623 .0497 .0421 .0367 .0300 .0275	5.0000 4.9887 4.9967 4.9930 5.0053 4.9933 5.0043 5.0026	.1481 .0915 .0622 .0507 .0433 .0366 .0323 .0288	

Histograms to summarize the distributions were made for each of the four estimators of variance. The next three plots show the histograms of the first approximation $v_1(q)$ for sample sizes 3, 5 and 9. The histograms for $v_2(q)$, mse₂(q) and $v(q^*)$ were similar. As the sample size increases, the distribution of the variance becomes less skewed. This behavior is similar to that of the gamma distribution for increasing values of the shape parameter.

A plotting technique similar to normal plotting has been designed for the gamma distribution [7]. This technique can be used to determine whether a random sample of observations come from a gamma distribution; if the distribution is gamma, the points will yield a straight line configuration.

Gamma plots were made of $v_1(q)$ for each sample size, and the last three plots for n = 3, 6 and 9, show that even for small sample sizes it is not unreasonable to assume that the distribution of the estimate of variance of q is chi-square.

The mean and variance of the sampling distributions of the four estimators of V(q) are given in Tables 2a and 2b.

TABLE 2a

Mean of 1000 Samples of the Estimates of V(q) for Small Sample Size: x Distribution Normal

n	v 1(d)	v ₂ (q)	mse(q)	v(q*)	V(y/x)
2 34 56 78 9	.1285	.1394	.1550	.1481	.1235
	.0853	.0892	.0927	.0915	.0896
	.0607	.0620	.0636	.0622	.0623
	.0490	.0498	.0509	.0507	.0497
	.0424	.0428	.0436	.0433	.0421
	.0354	.0356	.0362	.0366	.0367
	.0318	.0320	.0324	.0323	.0300
	.0283	.0284	.0324	.0288	.0275

TABLE 2b

Variance of 1000 Samples of the Estimates of V(q) for Small Sample Size: x Distribution Normal

n 	<u>v_l(q)</u>	v₂(q)	mse(q)	v(q *)
2 34 56 78 9	.0317	.0392	.0472	.0488
	.0072	.0080	.0088	.0096
	.0025	.0026	.0028	.0027
	.0014	.0014	.0015	.0017
	.0007	.0007	.0007	.0008
	.0004	.0004	.0004	.0005
	.0003	.0003	.0003	.0003
	.0002	.0002	.0002	.0002

Some improvement over $v_1(q)$ might be achieved by using $v_2(q)$ or $v(q^*)$. The mean of the one thousand $v_2(q)$'s is closer than $v_1(q)$ to the true variance for n = 3,4,5, and 7, the mean of the $v(q^*)$ is closer for n = 3,4 and 7. However, the precision of $v_1(q)$ is never less than the precision of any other estimator, and for n = 2,3 and 4, the precision of $v_1(q)$ is actually greater.

Thus, when y and x are normally distributed, $\overline{y}/\overline{x}$ and the first approximation $v_1(q)$, are good estimators of $\overline{Y}/\overline{X}$ and $V(\overline{y}/\overline{x})$ respectively. Any improvement in estimating the true variance by $v_2(q)$ or $v(q^*)$ will most likely not warrant In order to examine the behavior of the variance of q for non-normal populations, a second study was made similar to the first. Here, x was chi-square with 2 degrees of freedom, and a constant added such that E(x) = 12, V(x) = 4; y was defined as before so that E(y) = 60, V(y) = 125. The correlation between y and x was again equal to .89. For this population also q was an unbiased estimator of the population ratio.

Tables 3a and 3b give the mean and variance of the four estimators of V(q) for each sample size.

TABLE 3a

Mean of 1000 Samples of the Estimates of V(q) for Small Sample Size n: x Distribution Exponential

n	v _1(q)	v ₂ (q)	mse(q)	v (q*)	V(y/x)
2 34 56 78 9	.1096	.1134	.1143	.0728	.0734
	.0769	.0777	.0780	.0527	.0541
	.0561	.0562	.0563	.0393	.0399
	.0460	.0460	.0460	.0318	.0327
	.0402	.0401	.0402	.0281	.0277
	.0336	.0336	.0336	.0235	.0240
	.0306	.0305	.0305	.0214	.0204
	.0272	.0271	.0271	.0190	.0182

TABLE 3b

Variance of 1000 Samples of the Estimates of V(q) for Small Sample Size n: x Distribution Exponential

n 	<u>v₁(q)</u>	v ₂ (q)	mse(q)	v(q *)
2 34 56 78 9	.0245	.0263	.0268	.0084
	.0063	.0065	.0066	.0023
	.0023	.0023	.0023	.0009
	.0013	.0013	.0013	.0005
	.0006	.0006	.0006	.0003
	.0004	.0004	.0004	.00015
	.0003	.0003	.0003	.0001
	.0002	.0002	.0002	.00007

In this case the estimators $v_1(q)$, $v_2(q)$ and $mse_2(q)$ all lead to serious overestimates of the true variance. The means of these three estimators are about equal. The ratio of $v_1(q)$ to the true variance varies about 1.5 and the two estimators $v_2(q)$ and $mse_2(q)$ are actually worse than $v_1(q)$ for n = 2,3,4; for n = 8 and 9 It appears that they may begin to improve over $v_1(q)$. The precision of these three estimators is approximately the same. However, the mean of $v(q^*)$ is consistently near the true value of the variance and its precision is much greater than the other estimators.

The normal plots, gamma plots and histograms were similar in behavior to those of the first study.

The maximum observations of $v_1(q)$ and $v(q^*)$ are given in Table 4; the minimum values were approximately equal.

TABLE 4

Largest Observation From the Distributions of v₁(q) and v(q*) for Small Sample Sizes: x Distribution Exponential

n	Maximum v _l (q)	Maximum v(q*)
23456780	1.1249 .6689 .4256 .2358 .1772 .1284 .1128 .0012	.5282 .3385 .2000 .1254 .0922 .0697 .0670 .0526
7	•0912	·0520

Table 4 shows that the spread of $v_1(q)$ is nearly twice that of $v(q^*)$ but aside from this the distributions of the four estimators of the true variance are approximately the same, specifically, chi-square, even when the original distributions are badly skewed.

From the results of these two studies, it may be inferred that the bias of the estimator $v_1(q)$ is dependent upon the degree of skewness of the original y and x populations. Estimates of the true variance taken from higher order approximations lead only to slight improvements over the second order approximation $v_1(q)$, and in some cases the estimate is actually worse. The precision of $v(q^*)$ is nearly double that of $v_1(q)$ for exponential x distributions and the bias of $v(q^*)$ is smaller than that of $v_1(q)$. Thus it appears that the split-sample estimator q^* may be definitely preferable to q in some situations.

Computations in the Monte-Carlo study were done on the IBM 7090 computer at Bell Telephone Laboratories in Murray Hill. Plots were drawn by the Stromberg-Carlson 4020 microfilm printer using output from the 7090.

REFERENCES

 Cochran, W. G., <u>Sampling Techniques</u>, New York, John Wiley and Sons, Inc., 1953.

- [2] Goodman, L. A. and Hartley, H. O., "The Precision of Unbiased Ratio-Type Estimators," Journal of the American Statistical Association, 53 (1958), 491-508.
- [3] Kish, L., Namboodiri, N. J., and Pillai, R. K., "The Ratio Bias in Surveys," <u>Journal of the American Statistical Association</u>, 57 (1962), 863, 867.
- [4] Quenouille, M. H., "Notes on Bias in Estimation," <u>Biometrika</u>, 43 (1956) 353-360.
- [5] Sukhatme, P. V., Sample Theory of Surveys with Applications, Ames, <u>The</u> <u>Iowa State College Press</u>, 1954.
- [6] Tukey, J. W., "Bias and Confidence in Not-Quite Large Samples," <u>The</u> <u>Annals of Mathematical Statistics</u>, 29 (1958), 614.
- [7] Wilk, M. B., Gnanadesikan, R., and Huyett, Miss M. J., "Probability Plots for the Gamma Distribution," <u>Technometrics</u>, 4 (1962), 1-20.



Figure 2



-

Figure 4







Figure 6



Figure 7





Figure 9