

# Using Prediction-Oriented Software for Survey Estimation - Part II: Ratios of Totals

James R. Knaub, Jr.

US Dept. of Energy, Energy Information Administration, EI-53.1

**ABSTRACT:** This article is an extension of Knaub (1999), "Using Prediction-Oriented Software for Survey Estimation," which dealt with the estimation of totals and subtotals and the corresponding estimates of variance in the presence of 'missing data,' whether missing as part of a sampling scheme, or as a result of nonresponse. The current article deals with ratios of totals. An example from the electric power industry would be the estimation of revenue per kilowatthour and its associated variance estimate. As in Knaub (1999), the goal is to produce such estimates by making use of currently available software in which the model can be quickly and easily modified, and the data may be stored in such a manner that they may be easily recategorized for purposes of publishing various aggregations of the data with corresponding variance estimates. As in Knaub (1999), a blending of survey statistics and econometrics can be seen.

## **KEYWORDS:**

survey sampling; prediction; estimation; imputation; variance estimation; ratios of totals

## **SOME APPLICATIONS:**

A great advantage with this new method is that it is easy to store and manipulate data. Statistical agencies may present data in different formats, in different tables. It is cumbersome to aggregate data one way to estimate subtotals for one table and another way to estimate subtotals for overlapping areas for another table. The grand totals, for example, would differ. In the case of publishing subtotals for (Bureau of the) Census division regions, consider that Census divisions are groups of States. However, North American Electric Reliability Regions (NERC Regions) have boundaries that cut through States. Further, NERC boundaries recently moved. If imputed values were substituted for each 'missing' observation, using the largest, relatively homogeneous set of data available for each prediction, then they could be aggregated however desired, and using the method of Knaub(1999) and this article, standard errors may be estimated for any aggregation.

For establishment surveys, a very strong reason for using cutoff model sampling at all is that the smallest and most 'plentiful' establishments are often unable to supply data on a frequent basis with reasonable accuracy. A lot of imputation may be necessary. Resources are another problem. The method of this article, and Knaub(1999), also applies to imputation for census surveys, and may be used to help publish preliminary subtotals/totals and/or ratios of such numbers more timely. For a design-based sample, this method could be used to predict/impute for the missing members of the sample, and then the aggregate level variances for that part could be added to the variance estimates from the design-based sample. (This technique is apparently used elsewhere. See Lee, Rancourt, and Saerndal(1999).)

## **NEW METHODOLOGY:**

As shown in Knaub (1999), any statistical software package that will provide predicted values, a standard error or variance of the prediction error, and the mean square error (MSE) from the analysis of

variance, will suffice for estimating (sub)totals and their variances in the presence of ‘missing’ data, using the method found in that article. The regression weight must be supplied by means of considerations such as those found in Knaub (1997). For purposes of predicting missing numbers, the population should be categorized into the largest, relatively homogeneous sets of data possible. Imputed numbers are then each individually associated with variance related information that can be regrouped according to whatever aggregations one may wish to publish. The current article goes a step farther and associates pairs of numbers whose ratio is of interest, and then assigns covariance information to the pair for later aggregations. As in Knaub (1999), a given aggregation could contain little or no observed values, yet it may be possible to estimate totals or ratios of totals with some usefulness. Thus ‘small area statistics’ results may be available.

Here we consider  $V_L^*(T^* - T)$ , the variance of the error when estimating a total. This is a multiple regression form of  $V_L$  in Royall and Cumberland (1981), which contained some more robust variance estimates. However, Knaub (1992), page 879, Figure 1 shows that  $V_L$  may do very well, and this multiple regression form of this variance estimator has performed well, as in Knaub (1996) and Knaub (1999).

Now, according to Knaub (1999), using  $V_L^*(y_i^* - y_i)$  for the variance of the prediction error (see Maddala (1992)), and noting that  $V_L^*(T^* - T) = V_L^*(y_i^* - y_i)$  when there is only one missing value, one finds in Knaub (1999) that in general, we may approximate as follows:

$$V_L^*(T^* - T) = \delta(N - n) \sum_r \left\{ V_L^*(y_i^* - y_i) - \frac{\sigma_e^{*2}}{w_i} \right\} + \sum_r \frac{\sigma_e^{*2}}{w_i},$$

where,  $0 < \delta < 1$  ( $\delta = 0.3$  may be a fair general use value;

further discussion is found in Knaub (1999)),

$\sum_r$  means to sum over the cases with missing data (see Royall (1970)),

$\sigma_e^{*2}$  (Knaub (1996)) is the estimated variance of the random factor of the residual,  $e_0$ ,

(see Knaub (1993, 1995)), where the error term is  $e_i = w_i^{-1/2} e_{o_i}$ ,

$w_i$  is the regression weight, and

$(N - n)$  is the number of members of the population that are not in the sample.

(Note: As  $(N - n)$  approaches 1,  $\delta$  approaches 1. However,  $\delta$  will generally decrease quickly as  $(N - n)$  becomes a little larger.)

**Following is an excerpt from Knaub (1999), page 8:**

Picture a typical data file as follows, where "EG" is a category for purposes of performing predictions (an "estimation group"), and "PG" is a category for purposes of publishing subtotals (a "publication group"). Each line represents a record for a given member of the population. A  $y$  value is an observed (or "collected") value,

and  $y^*$  is a predicted value. Let  $S1_i^2 = \sum_L (y_i^* - y_i)$ , the variance of the prediction error, and

$S2_i^2 = \sigma_e^{*2} / w_i$ , the mean square error divided by the regression weight, for each case,  $i$ .

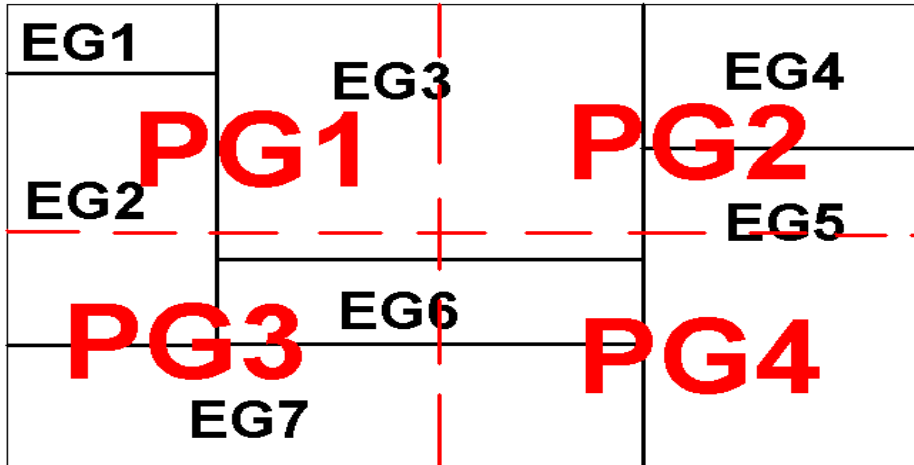
**Example of a partial file:**

$y_i$ or $y_i^*$	$S1_i$	$S2_i$	EG	PG (a)	PG (b)	PG (c)
6725	0	0	1	1	2	5
4359	0	0	1	2	1	3
1289	0	0	2	1	4	4
497	20	17	1	1	3	2
317	13	11	1	2	2	2
278	10	9	1	1	3	2
223	9	8	2	1	3	2

After that, Knaub (1999) discusses adjustments for nonsampling error that would be applicable here also, but will not be repeated here.

A figure from Knaub (1999), page 12, is reprinted next, as it may help in the general visualization of the implementation of this method.

Consider, for example, that a population could be subdivided into seven categories which each represent nearly homogeneous data under one model per category. Therefore, there are seven EGs. Further suppose that the population is divided into four parts for which subtotals and their variances are to be estimated. That would mean having four PGs. This could be represented visually as in the figure below:



.....  
 The dashed lines represent the separation of the universe into the four "publication groups" ("PGs"). The solid lines represent boundaries between "estimation groups" ("EGs").

Here,  $y^*$ ,  $V_L^*(y_i^* - y_i)$ , and  $\sigma_e^{*2}/w_i$  are estimated for each missing observation within a given "EG" group, using all data in that group. Then every piece of an EG within a given PG is treated as a stratum for estimating the total for that PG group. The variance for each stratum is estimated using the  $V_L^*(T^* - T)$  formula above, and the total variance estimate is found by adding the strata variance

estimates. Let  $T^* = \sum_S y_i + \sum_R y_i^*$  be the total for any given stratum (*i.e.*, for a given 'piece' of an "EG" within a "PG").  $\sum_S$  refers to summation over the observed data, and  $\sum_R$  refers to summation over predictions for missing data (see Royall (1970)). The corresponding variance for that stratum is  $V_L^*(T^* - T) = \delta(N - n) \sum_R \left\{ V_L^*(y_i^* - y_i) - \frac{\sigma_e^{*2}}{w_i} \right\} + \sum_R \frac{\sigma_e^{*2}}{w_i}$ .

The current problem, however, is to extend this to estimating variance for the estimated ratios of such (sub)totals. Let the estimated ratio be  $T_A^* / T_B^*$ , and the variance sought is  $V_L^*(T_A^* / T_B^*)$ . In the case of totals, estimates of subtotals and their variances within strata are simply added to obtain an estimate of a total and its variance, respectively. In the case of an estimated ratio of totals, the numerator and denominator are estimated separately, adding stratum components until the estimates of numerator and denominator are completed, and then the estimated ratio is found. For the estimated variance of this estimated ratio, an estimated variance for the numerator, and an estimated variance for the denominator, and a covariance estimate will each have to be constructed from the strata estimates, and then applied to the overall estimations of the ratio and its variance. To do this, however, in a manner that is flexible to changes in PG categorizations,  $V_L^*(y_i^* - y_i)$ , and  $\sigma_e^{*2} / w_i$  will be needed for the numerator and for the denominator, and a fifth number, a covariance component, will be needed, for each related pair of missing data points in the population. This is very little data to have to store and yet leave such flexibility in the publication process for data aggregations.

For one stratum, the estimation of the ratio,  $T_A / T_B$ , designated  $T_A^* / T_B^*$ , is straightforward. We simply use  $T^* = \sum_S y_i + \sum_R y_i^*$  (Royall (1970)) for the numerator, and repeat the application for the denominator. However, variance estimation is more involved and will be discussed below. Further, when considering more than one stratum, an estimated total can be found by adding the subtotal estimates by strata, and similarly for the estimated variance of the total. The estimation of the ratio of totals is also straightforward, but would now involve a more complicated variance formula. However, variance estimation would still rely on only the five stored bits of information for every pair of missing data points.

## VARIANCE ESTIMATION:

Starting with the case of a single stratum:

Knaub (1994) is largely a review of and relies heavily upon P.S.R.S. Rao (1992) for covariance formulae associated with the variance of a ratio of variables. Here, however, as in Knaub (1999), the thrust is somewhat different. Here the emphasis is on simplicity of operation, including easily revised models and the association of all information at the individual (pairs of) point(s) level that will be needed to estimate ratios and their variances at any level of aggregation.

As in Knaub (1994),

$$\frac{V_L^*(T_A^*/T_B^*)}{(T_A^*/T_B^*)^2} = \frac{V_L^*(T_A^*)}{T_A^{*2}} + \frac{V_L^*(T_B^*)}{T_B^{*2}} - 2 \frac{\text{COV}_L^*(T_A^*, T_B^*)}{T_A^* T_B^*}$$

So,

$$V_L^*(T_A^*/T_B^*) = \frac{V_L^*(T_A^*)}{T_B^{*2}} + \frac{T_A^{*2} V_L^*(T_B^*)}{T_B^{*4}} - 2 \frac{T_A^* \text{COV}_L^*(T_A^*, T_B^*)}{T_B^{*3}}$$

Also from Knaub (1994) and Hansen, Hurwitz and Madow (1953), pages 56 to 58,

$$\text{COV}_L^*(T_A^*, T_B^*) = \sum_r \text{COV}_L^*(y_{Ai}^*, y_{Bi}^*) + \dots$$

which corresponds to

$$V_L^*(T^* - T) > \sum_r V_L^*(y_i^* - y_i) \quad \text{in Knaub (1999).}$$

Also, by Knaub (1999)

$$V_L^*(T^* - T) = \sum_r \sigma_e^{*2} / w_i + (N-n)^2 V^*(b_0) + \left( \sum_r x_{1i} \right)^2 V^*(b_1) + \dots \quad \text{and}$$

$$\sum_r V_L^* (y_i^* - y_i) = \sum_r \sigma_e^{*2} / w_i + (N-n) V^* (b_0) + \left( \sum_r x_i^2 \right) V^* (b_1) + \dots$$

By Knaub (1996),  $\sigma_e^{*2} = \sum_{i=1}^n e_{0i}^2 / \text{d.f.} = \sum_{i=1}^n w_i e_i^2 / \text{d.f.}$ ,

where  $e_i = y_i - y_i^* =$  residual and d.f. is the number of degrees of freedom, so

$$\text{COV}_L^* (y_{Ai}^*, y_{Bi}^*) = \frac{\sigma_{e; y_A, y_B}^*}{w_{Aj}^{0.5} w_{Bj}^{0.5}} + \dots = \frac{\sum w_{Aj}^{0.5} e_{Aj} w_{Bj}^{0.5} e_{Bj}}{w_{Aj}^{0.5} w_{Bj}^{0.5} (\text{d.f.})} + \dots$$

therefore

$$\text{COV}_L^* (T_A^*, T_B^*) \approx \left[ \sum_r \frac{\sigma_{e; y_A, y_B}^*}{w_{Ai}^{0.5} w_{Bi}^{0.5}} \right] \left[ \frac{V_L^* (T_A^*) \cdot V_L^* (T_B^*)}{\sum_r \frac{\sigma_{Ae}^{*2}}{w_{Ai}} \cdot \sum_r \frac{\sigma_{Be}^{*2}}{w_{Bi}}} \right]^{1/2}$$

or

$$\text{COV}_L^* (T_A^*, T_B^*) \approx \left[ \sum_r w_{Aj}^{-0.5} w_{Bj}^{-0.5} \right] \left[ \sum w_{Aj}^{0.5} e_{Aj} w_{Bj}^{0.5} e_{Bj} / \text{d.f.} \right] \left[ \frac{V_L^* (T_A^*) \cdot V_L^* (T_B^*)}{\sum_r \frac{\sigma_{Ae}^{*2}}{w_{Ai}} \cdot \sum_r \frac{\sigma_{Be}^{*2}}{w_{Bi}}} \right]^{1/2}$$

(Note: for  $e_{Ai}$ , and  $e_{Bi}$ , one can save  $y_i - y_i^*$  in each case (A and B) in another file.)

So, in addition to producing  $V_L^*(y_i^* - y_i)$ , and  $\sigma_e^{*2}/w_i$ , for each ‘missing’ number, also save

$$\frac{\sigma_{e; y_A, y_B}^*}{w_{Ai}^{0.5} w_{Bi}^{0.5}} \text{ for each pair of corresponding, missing numbers.}$$

Once again, per Knaub (1994) and Hansen, Hurwitz and Madow (1953), sum over  $COV_{L_k}^*(T_{Ak}^*, T_{Bk}^*)$  just as is done for  $V_{L_k}^*(T_k^* - T_k)$  for the case of multiple strata,  $k$ .

$$\text{Then } V_L^*(T_A^*/T_B^*) = \frac{V_L^*(T_A^*)}{T_B^{*2}} + \frac{T_A^{*2} V_L^*(T_B^*)}{T_B^{*4}} - 2 \frac{T_A^* COV_L^*(T_A^*, T_B^*)}{T_B^{*3}}, \text{ where}$$

$$T_A^* = \sum_k T_{Ak}^*, \quad T_B^* = \sum_k T_{Bk}^*,$$

$$V_L^*(T_A^*) = \sum_k V_{L_k}^*(T_{Ak}^* - T_{Ak}), \quad V_L^*(T_B^*) = \sum_k V_{L_k}^*(T_{Bk}^* - T_{Bk}) \text{ and}$$

$$COV_L^*(T_A^*, T_B^*) = \sum_k COV_{L_k}^*(T_{Ak}^*, T_{Bk}^*), \text{ and remembering that}$$

$$T_A^* = \sum_S y_{Ai} + \sum_r y_{Ai}^*, \quad T_B^* = \sum_S y_{Bi} + \sum_r y_{Bi}^*,$$

$$V_{L_k}^*(T_{Ak}^* - T_{Ak}) = \delta_A (N_A - n_A) \sum_r \left\{ V_{L_k}^*(y_{Ak_i}^* - y_{Ak_i}) - \frac{\sigma_{Ae}^{*2}}{w_{Ai}} \right\} + \sum_r \frac{\sigma_{Ae}^{*2}}{w_{Ai}},$$

$$V_{L_k}^*(T_{Bk}^* - T_{Bk}) = \delta_B (N_B - n_B) \sum_r \left\{ V_{L_k}^*(y_{Bk_i}^* - y_{Bk_i}) - \frac{\sigma_{Be}^{*2}}{w_{Bi}} \right\} + \sum_r \frac{\sigma_{Be}^{*2}}{w_{Bi}},$$

and



$$\text{COV}_L^*(T_A^*, T_B^*) \approx \left[ \sum_r w_{Aj}^{-0.5} w_{Bj}^{-0.5} \right] \left[ \sum^n w_{Aj}^{0.5} e_{Aj} w_{Bj}^{0.5} e_{Bj} / \text{d.f.} \right] \left[ \frac{V_L^*(T_A^*)}{\sum_r \frac{\sigma_{Ae}^{*2}}{w_{Ai}}} \cdot \frac{V_L^*(T_B^*)}{\sum_r \frac{\sigma_{Be}^{*2}}{w_{Bi}}} \right]^{1/2} .$$

So, using  $\frac{\sigma_{e; y_A, y_B}^*}{w_{Ai}^{0.5} w_{Bi}^{0.5}} = S3_i^2$ , an example of the requisite data file is as follows:

**Example of a partial file:**

$y_{Ai}$	$y_{Bi}$		$y_{Bi}$	$y_{Bi}$			EG	PG(a)	PG(b)
or $y_{Ai}^*$	$S1_{Ai}$	$S2_{Ai}$	or $y_{Bi}^*$	$S1_{Bi}$	$S2_{Bi}$	$S3_i$			
6725	0	0	432	0	0	0	1	1	2
4359	0	0	320	0	0	0	1	2	1
1289	0	0	85	0	0	0	2	1	4
497	20	17	35	9	8	3	1	1	3
317	13	11	22	6	5	2	1	2	2
278	10	9	17	3	3	1	1	1	3
223	9	8	14	3	2	1	2	1	3

**Example application:**

Electricity sales and associated revenue data for the residential economic end-use sector were taken from a census survey of utilities in Indiana and Ohio for two succeeding years in which the survey was conducted. For the more recent year, according to these data (which would, naturally, include some nonsampling error), the overall revenue per kilowatthour found for Indiana was 7.010 cents per kilowatthour, and for Ohio it was 8.704 cents per kilowatthour, for the residential sector.

There were 240 utilities in these two States which were used for this experiment. They each fit under one of three ownership classes (municipally owned, privately owned or cooperatives), and modeling by

ownership code showed substantially different growth rates in revenue, depending upon ownership class. (Sales probably also would show this difference clearly, but multiple regressors were used for sales, preventing direct comparison of the coefficients for a single regressor, but exercising the flexibility of the methodology.) Therefore, there were three estimation groups (EG1, EG2, EG3), the ownership classes, and two publication groups (PG1, PG2), which were the two States. Data from the most recent census were prepared as follows: Every fourth response was removed as if there were a nonresponse rate of 25%. Four such data sets were prepared, similarly, but with a different start (first, second, third or fourth record) for the process of designating which numbers to eliminate as if they were nonresponses. Therefore, there was one data set used for regressor data, and four others for the data of interest for four tests. Between the four data sets to be used for the ‘current’ data, every observation was eliminated once and only once. Therefore, a particularly favorable data set (with respect to test results) must be balanced by one or more less favorable data sets. The following tables show some results of this experiment. What is of greatest interest here is the difference between standard error estimates, with and without a nonzero covariance term. Obtaining an estimate when covariance should be considered nonzero was the point of this extension to Knaub (1999).

These data are highly skewed. Totals for sales and revenue in Ohio are dominated by eight very large utilities. Thus, a 25% nonresponse could sometimes mean that the majority of an estimated total is imputed.

## Indiana

	cents per kilowatthour from full test data set	estimated cents per kilowatthour	‘error’	estimated standard error	estimated standard error if zero covariance
Test 1	7.010	6.984	0.026	0.022	0.080
Test 2	7.010	7.020	-0.010	0.021	0.044
Test 3	7.010	7.011	-0.001	0.025	0.060
Test 4	7.010	7.011	-0.001	0.003	0.020

## Ohio

	cents per kilowatthour from full test data set	estimated cents per kilowatthour	‘error’	estimated standard error	estimated standard error if zero covariance
Test 1	8.704	8.695	0.009	0.002	0.004
Test 2	8.704	8.736	-0.032	0.043	0.073
Test 3	8.704	8.706	-0.002	0.004	0.012
Test 4	8.704	8.698	0.006	0.036	0.080

As previously stated, the test data came from a set of 240 observations of electric utility sales and revenue data (residential) in Indiana and Ohio. State and ownership ‘type’ were as follows:

	<b>Municipal</b>	<b>Private</b>	<b>Cooperative</b>	
<b>Indiana</b>	42	6	73	121
<b>Ohio</b>	27	8	84	119
	69	14	157	<b>240</b>

Every fourth one of the 240 observations was treated as a nonresponse. The four complementary data test sets were prepared by using a different ‘start’ record (i.e., first, second, third or fourth) for each.

**Software:**

Computer code was shown and discussed in Knaub (1999), but the situation is somewhat more complex here. More manipulation may be needed to use residuals in accordance with the equations shown in this article, so that covariance is taken into account. It is still likely, however, to be relatively easy to make use of existing software, as was done for the example above. Further, the ease with which data may be stored and used for multiple purposes by post-production software makes this method attractive.

If your software calculates mean square error, as shown in the code on page 34 in Knaub(1999), this is no longer adequate when estimating covariance. Each residual needs to be identified. (See the covariance formula at the bottom of page 7 in the current article.)

If regression weights are passed from one computer program to another, be certain to save enough significant digits. It may be best to save the square root of these weights. (This is usually not a problem in survey statistics where the problem more often is the presentation of too many digits for a given statistic.)

For more information, please contact the author.

**Generalized  $\delta$  :**

Values for delta were also discussed in Knaub(1999). When N-n is small, the value of delta should be increased, but in such cases, the impact on variance from the term where delta is found is correspondingly decreased. Considering delta to be a constant, as in Knaub(1999), say  $\delta = q$ , is probably adequate, but the following could be used:  $\delta = q + (1 - q^2)(1 + q)^{-(N-n)}$ .

## REFERENCES:

- Hansen, M.H., Hurwitz, W.N., and Madow, W.G. (1953), Sample Survey Methods and Theory, Volume II: Theory, John Wiley & Sons.
- Knaub, J.R., Jr. (1992), "More Model Sampling and Analyses Applied to Electric Power Data," Proceedings of the Section on Survey Research Methods, American Statistical Association, pp. 876-881.
- Knaub, J.R., Jr. (1993), "Alternative to the Iterated Reweighted Least Squares Method: Apparent Heteroscedasticity and Linear Regression Model Sampling," Proceedings of the International Conference on Establishment Surveys, American Statistical Association, pp. 520-525.
- Knaub, J.R., Jr. (1994), "Relative Standard Error for a Ratio of Variables at an Aggregate Level Under Model Sampling," Proceedings of the Section on Survey Research Methods, Vol. I, American Statistical Association, pp. 310-312.
- Knaub, J.R., Jr. (1995), "A New Look at 'Portability' for Survey Model Sampling and Imputation," Proceedings of the Section on Survey Research Methods, Vol. II, American Statistical Association, pp. 701-705.
- Knaub, J.R., Jr. (1996), "Weighted Multiple Regression Estimation for Survey Model Sampling," InterStat, May 1996, <http://interstat.stat.vt.edu>. (Note shorter, more recent version in ASA Survey Research Methods Section proceedings, 1996.)
- Knaub, J.R., Jr. (1997), "Weighting in Regression for Use in Survey Methodology," InterStat, April 1997, <http://interstat.stat.vt.edu>. (Note shorter, more recent version in ASA Survey Research Methods Section proceedings, 1997.)
- Knaub, J.R., Jr. (1999), "Using Prediction-Oriented Software for Survey Estimation," InterStat, August 1999, <http://interstat.stat.vt.edu>, partially covered in "Using Prediction-Oriented Software for Model-Based and Small Area Estimation," to appear in ASA Survey Research Methods Section proceedings, 1999.
- Lee, H., Rancourt, E., and Saerndal, C.-E. (1999), "Variance Estimation from Survey Data Under Single Value Imputation," presented at the International Conference on Survey Nonresponse, Oct. 1999, to be published in a monograph.
- Maddala, G.S. (1992), Introduction to Econometrics, 2<sup>nd</sup> ed., Macmillan Pub. Co.
- Rao, Poduri S.R.S. (1992), unpublished letters, Aug. - Oct. 1992, on covariances associated with three Royall and Cumberland model sampling variance estimators.
- Royall, R.M. (1970), "On Finite Population Sampling Theory Under Certain Linear Regression Models," Biometrika, 57, pp. 377-387.
- Royall, R.M. and Cumberland, W.G. (1981), "An Empirical Study of the Ratio Estimator and Estimators of its Variance," Journal of the American Statistical Association, 76, pp.66-88.